



International

Rapport sur la sécurité de l'IA

Le rapport scientifique international
sur la sécurité de l'IA avancée

Janvier 2025

Contributeurs

CHAISE

Prof. Yoshua Bengio, Université de Montréal / Mila - Institut québécois d'IA

GROUPE CONSULTATIF D'EXPERTS

Ce panel international a été nommé par les gouvernements des 30 pays énumérés ci-dessous, l'ONU, l'UE et l'OCDE.

Australie : Bronwyn Fox, Université de New

Sud du Pays de Galles

Brésil : André Carlos Ponce de Leon Ferreira de

Carvalho, Institut de mathématiques et

Informatique, Université de São Paulo

Canada : Mona Nemer, conseillère scientifique en chef

du Canada

Chili : Raquel Pezoa Rivera, Universidad

Technicien Federico Santa Maria

Chine : Yi Zeng, Académie chinoise des sciences

Union européenne : Juha Heikkilä, IA européenne

Bureau

France : Guillaume Avrin, Coordination nationale

pour l'intelligence artificielle

Allemagne : Antonio Krüger, Recherche allemande

Centre d'intelligence artificielle

Inde : Balaraman Ravindran, École Wadhvani de science

des données et d'IA, Institut indien de

Technologie Madras

Indonésie : Hammam Riza, collaboratif

Recherche et innovation industrielle en IA

Renseignement (KORIKA)

Irlande : Ciarán Seoighe, Research Ireland

Israël : Ziv Katzir, Autorité israélienne de l'innovation

Italie : Andrea Monti, expert juridique pour la

Sous-secrétaire d'État au numérique

Transformation, Conseil des ministres italiens

Présidence

Japon : Hiroaki Kitano, Sony Group Corporation

Kenya : Nusu Mwamanzi, ministère des TIC et

Économie numérique

Royaume d'Arabie Saoudite : Fahad Albalawi,

Autorité saoudienne des données et de l'intelligence artificielle

Intelligence

Mexique : José Ramón López Portillo, LobsterTel

Pays-Bas : Haroon Sheikh, Pays-Bas

Conseil scientifique pour la politique gouvernementale

Nouvelle-Zélande : Gill Jolly, ministère des Affaires,

Innovation et emploi

Nigéria : Olubunmi Ajala, Ministère de

Communication, Innovation et Digital

Économie

OCDE : Jerry Sheehan, directeur de l'

Direction de la science, de la technologie et

Innovation

Philippines : Dominic Vincent Ligot, CirroLytix

République de Corée : Kyoung Mu Lee,

Département d'électricité et d'informatique

Ingénierie, Université nationale de Séoul

Rwanda : Crystal Rugege, Centre pour la Quatrième Révolution industrielle

Singapour : Denise Wong, Data Innovation et Groupe de protection, Infocomm Media Autorité de développement

Espagne : Nuria Oliver, ELLIS Alicante

Suisse : Christian Busch, conseiller fédéral Département des affaires économiques, de l'éducation et de la recherche

Turquie : Ahmet Halit Hatip, ministère turc de Industrie et technologie

Ukraine : Oleksii Molchanovskyi, expert Comité pour le développement de l'intelligence artificielle Renseignements en Ukraine

Émirats arabes unis : Marwan Alserkal, Ministère des Affaires du Cabinet, Cabinet du Premier ministre Bureau

Royaume-Uni : Chris Johnson, chef Conseiller scientifique au Département de Science, innovation et technologie

Nations Unies : Amandeep Singh Gill, Secrétaire général adjoint chargé du numérique et Technologies émergentes et Envoyé du Secrétaire général pour la technologie

États-Unis : Saif M. Khan, ministère américain du Commerce

DIRECTEUR SCIENTIFIQUE

Sören Mindermann, Mila - Institut québécois d'IA

RÉDACTEUR PRINCIPAL

Daniel Privitera, Centre KIRA

GROUPE D'ÉCRITURE

Tamay Besiroglu, Epoch AI

Rishi Bommasani, Université de Stanford

Stephen Casper, Institut du Massachusetts Technologie

Yejin Choi, Université de Stanford

Philip Fox, Centre KIRA

Ben Garfinkel, Université d'Oxford

Danielle Goldfarb, Mila - Institut québécois d'IA

Hoda Heidari, Université Carnegie Mellon

Anson Ho, Epoch AI

Sayash Kapoor, Université de Princeton

Leila Khalatbari, Université de Hong Kong Science et technologie

Shayne Longpre, Institut du Massachusetts Technologie

Sam Manning, Centre pour la gouvernance de l'IA

Vasilios Mavroudis, Institut Alan Turing

Mantas Mazeika, Université de l'Illinois à Urbana-Champaign

Julian Michael, Université de New York

Jessica Newman, Université de Californie, Berkeley

Kwan Yee Ng, Concordia AI

Chinasa T. Okolo, Brookings Institution

Deborah Raji, Université de Californie, Berkeley

Girish Sastry, indépendant

Elizabeth Seger (écrivain généraliste), Demos

Theodora Skeadas, Intelligence humaine

Tobin South, Institut du Massachusetts

Technologie

Emma Strubell, Université Carnegie Mellon

Florian Tramèr, ETH Zurich

Lucia Velasco, Université de Maastricht

Nicole Wheeler, Université de Birmingham

CONSEILLERS PRINCIPAUX

Daron Acemoglu, Institut du Massachusetts

Technologie

Olubayo Adekanmbi, a contribué en tant que Senior
Conseiller avant de prendre ses fonctions chez EqualzAI

David Dalrymple, Recherche avancée +
Agence d'invention

Thomas G. Dietterich, Université d'État de l'Oregon

Edward W. Felten, Université de Princeton

Pascale Fung, a contribué en tant que conseillère principale
avant de prendre ses fonctions chez Meta

Pierre-Olivier Gourinchas, Recherche
Département, Fonds monétaire international

Fredrik Heintz, Université de Linköping

Geoffrey Hinton, Université de Toronto

Nick Jennings, Université de Loughborough

Andreas Krause, ETH Zurich

Susan Leavy, Collège universitaire de Dublin

Percy Liang, Université de Stanford

Teresa Ludermir, Université fédérale de
Pernambouc

Vidushi Marda, collaboration IA

Helen Margetts, Université d'Oxford

John McDermid, Université de York

Jane Munga, Fondation Carnegie pour
Paix internationale

Arvind Narayanan, Université de Princeton

Alondra Nelson, Institut d'études avancées

Clara Neppel, IEEE

Alice Oh, École d'informatique de KAIST

Gopal Ramchurn, Responsable AI Royaume-Uni

Stuart Russell, Université de Californie,
Berkeley

Marietje Schaake, Université de Stanford

Bernhard Schölkopf, Institut ELLIS de Tübingen

Dawn Song, Université de Californie, Berkeley

Alvaro Soto, Université Pontificale Catholique de
Chili

Lee Tiedrich, Université Duke

Gaël Varoquaux, Inria

Andrew Yao, Institut de recherche interdisciplinaire
Sciences de l'information, Université Tsinghua

Ya-Qin Zhang, Université Tsinghua

SECRÉTARIAT

Institut de sécurité de l'IA

Baran Acar

Ben Clifford

Lambrini Das

Claire Dennis

Freya Hempleman

Hannah Marchand

Rian Overy

Ben Snodin

Mila — Institut québécois d'IA

Jonathan Barry

Benjamin Prud'homme

REMERCIEMENTS

Examineurs de la société civile et de l'industrie

Société civile : Institut Ada Lovelace, AI Forum New Zealand / Te Kāhui Atamai lahiko o Aotearoa, Groupe temporaire d'experts en IA d'Australie, Carnegie Endowment for International Peace, Center for Law and Innovation / Certa Foundation, Centre pour la gouvernance de l'IA, Chief Justice Meir Shamgar Center for Digital Law and Innovation, Institut Eon, Gradient Institute, Israel Democracy Institute, Fondation Mozilla, Old Ways New, RAND, SaferAI, Centre pour la résilience à long terme, The Future Society, Institut Alan Turing, Royal Society, Association turque des politiques d'intelligence artificielle.

Secteur d'activité : Advai, Anthropic, Cohere, Deloitte Consulting USA et Deloitte LLM UK, G42, Google DeepMind, Harmony Intelligence, Hugging Face, IBM, Lelapa AI, Meta, Microsoft, Shutterstock, Zhipu.ai.

Remerciements spéciaux

Le Secrétariat apprécie le soutien, les commentaires et les retours d'information d'Angie Abdilla, Concordia AI, Nitarshan Rajkumar, Geoffrey Irving, Shannon Vallor, Rebecca Finlay et Andrew Strait.

© Propriété de la Couronne 2025

Cette publication est sous licence selon les termes de la Licence gouvernementale ouverte v3.0, sauf indication contraire. Pour consulter cette licence, visitez <https://www.nationalarchives.gov.uk/doc/open-government-licence/version/3/> ou écrivez à l'équipe de politique d'information, Archives nationales, Kew, Londres TW9 4DU, ou envoyez un e-mail à : psi@nationalarchives.gsi.gov.uk.

Si nous avons identifié des informations relatives aux droits d'auteur appartenant à des tiers, vous devrez obtenir l'autorisation des détenteurs des droits d'auteur concernés.

Toute demande de renseignements concernant cette publication doit être envoyée à :

Secretariat.AIStateofScience@dsit.gov.uk

Les demandes de renseignements concernant le contenu du rapport doivent également être adressées au [responsable scientifique](#).

Clause de non-responsabilité

Le présent rapport ne reflète pas les opinions du président, d'un membre particulier du groupe de rédaction ou du groupe consultatif, ni d'aucun des gouvernements qui ont soutenu son élaboration. Il s'agit d'une synthèse des recherches existantes sur les capacités et les risques de l'IA avancée. Le président du rapport en est le responsable ultime et a supervisé son élaboration du début à la fin.

Numéro de série de recherche : DSIT 2025/001

Avant-propos	8
À propos de ce rapport	10
Point sur les dernières avancées de l'IA après la rédaction de ce rapport : note du président	11
Principales conclusions du rapport	13
Résumé exécutif	15
Introduction	25
Capacités de l'IA à usage général	29
1.1. Comment l'IA à usage général est-elle développée	30
1.2. Capacités actuelles	37
1.3. Capacités dans les années à venir	46
Risques	61
2.1. Risques liés à une utilisation malveillante	62
2.1.1. Atteintes aux personnes par le biais de faux contenus	62
2.1.2. Manipulation de l'opinion publique	67
2.1.3. Cyberinfraction	72
2.1.4. Attaques biologiques et chimiques	79
2.2. Risques liés aux dysfonctionnements	88
2.2.1. Problèmes de fiabilité	88
2.2.2. Biais	92
2.2.3. Perte de contrôle	100
2.3. Risques systémiques	110
2.3.1. Risques liés au marché du travail	110
2.3.2. Fracture mondiale de la recherche et du développement en IA	119
2.3.3. Concentration du marché et points de défaillance uniques	123
2.3.4. Risques pour l'environnement	128
2.3.5. Risques pour la vie privée	139
2.3.6. Risques de violation du droit d'auteur	144
2.4. Impact des modèles d'IA polyvalents à pondération ouverte sur les risques liés à l'IA	149
Approches techniques de la gestion des risques	157
3.1. Aperçu de la gestion des risques	158
3.2. Défis généraux en matière de gestion des risques et d'élaboration des politiques	169
3.2.1. Défis techniques pour la gestion des risques et l'élaboration des politiques	169
3.2.2. Défis sociétaux pour la gestion des risques et l'élaboration des politiques	176
3.3. Identification et évaluation des risques	181
3.4. Atténuation et surveillance des risques	191
3.4.1. Former des modèles plus fiables	191
3.4.2. Suivi et intervention	201
3.4.3. Méthodes techniques de protection de la vie privée	208
Conclusion	214
Liste des acronymes	216
Glossaire	218
Comment citer ce rapport	229
Références	230



Professeur Yoshua Bengio
Université de Montréal / Mila –
Institut et Chaire d'IA du Québec

Construire une compréhension scientifique partagée dans un domaine en évolution rapide

J'ai l'honneur de présenter le rapport international sur la sécurité de l'IA. Il s'agit du travail de 96 experts internationaux en IA qui ont collaboré dans un effort sans précédent pour établir une compréhension scientifique partagée à l'échelle internationale des risques liés à l'IA avancée et des méthodes de gestion. eux.

Nous avons entrepris ce voyage il y a un peu plus d'un an, peu après que les pays présents au Sommet sur la sécurité de l'IA de Bletchley Park aient accepté de soutenir la création de ce rapport. Depuis lors, nous avons publié un rapport intermédiaire en mai 2024, qui a été présenté au Sommet de Séoul sur l'IA. Nous sommes maintenant heureux de publier le présent rapport complet en prévision du Sommet sur l'action en matière d'IA qui se tiendra à Paris en février 2025.

Depuis le sommet de Bletchley, les capacités de l'IA à usage général, le type d'IA sur lequel se concentre ce rapport, ont encore augmenté. Par exemple, de nouveaux modèles ont montré des performances nettement meilleures lors de tests de programmation et de raisonnement scientifique. En outre, de nombreuses entreprises investissent désormais dans le développement d'« agents » d'IA à usage général.

des systèmes capables de planifier et d'agir de manière autonome pour atteindre des objectifs avec peu ou pas de surveillance humaine.

S'appuyant sur le rapport intermédiaire (mai 2024), le présent rapport reflète ces nouveaux développements. En outre, les experts ayant contribué à ce rapport ont apporté plusieurs autres modifications par rapport au rapport intermédiaire.

Par exemple, ils ont travaillé à améliorer davantage la rigueur scientifique de toutes les sections, ont ajouté une discussion sur des sujets supplémentaires tels que les modèles de pondération ouverte et ont restructuré le rapport pour qu'il soit plus pertinent pour les décideurs politiques, notamment en soulignant les lacunes en matière de preuves et les principaux défis pour les décideurs politiques.

Je tiens à exprimer ma profonde gratitude à l'équipe d'experts qui a contribué à ce rapport, notamment à nos rédacteurs, à nos conseillers principaux et au groupe consultatif international d'experts. J'ai été impressionné par leur excellence scientifique et leur expertise, ainsi que par l'attitude collaborative avec laquelle ils ont abordé ce projet ambitieux. Je suis également reconnaissant aux organisations de l'industrie et de la société civile qui ont examiné le rapport, en apportant des commentaires inestimables qui ont permis à ce rapport d'être plus complet qu'il ne l'aurait été autrement.

Je remercie également le gouvernement britannique d'avoir lancé ce processus et d'avoir offert un soutien opérationnel exceptionnel. Il était également important pour moi que le gouvernement britannique accepte que les scientifiques qui rédigent ce rapport bénéficient d'une totale indépendance.

L'IA reste un domaine en constante évolution. Pour suivre ce rythme, les décideurs politiques et les gouvernements doivent avoir accès aux connaissances scientifiques actuelles sur les risques que l'IA avancée pourrait présenter. J'espère que ce rapport ainsi que les publications futures aideront les décideurs à garantir que les citoyens du monde entier puissent bénéficier des avantages de l'IA en toute sécurité.

Pour tirer parti des opportunités de l'IA en toute sécurité, il faut une collaboration mondiale

Depuis la publication de la version intermédiaire de ce rapport, les capacités de l'IA avancée n'ont cessé de croître. Nous savons que cette technologie, si elle est développée et utilisée de manière sûre et responsable, offre des opportunités extraordinaires : développer nos économies, moderniser nos services publics et améliorer la vie de nos citoyens. Pour saisir ces opportunités, il est impératif que nous approfondissions notre compréhension collective de la manière dont l'IA peut être développée en toute sécurité.

Ce rapport historique témoigne de l'importance de la coopération internationale pour forger cette compréhension commune. Il est le fruit du travail de plus de 90 experts en IA de différents continents, secteurs et domaines d'expertise, qui se sont réunis pour offrir aux dirigeants et aux décideurs un point de référence mondial et un outil pour éclairer les politiques sur la sécurité de l'IA. Notre compréhension collective des systèmes d'IA de pointe s'est améliorée. Cependant, ce rapport souligne que l'IA de pointe reste un domaine de recherche scientifique actif, les experts continuant de diverger sur sa trajectoire et l'ampleur de son impact.

Nous maintiendrons l'élan de cet effort collectif pour faire avancer consensus scientifique mondial. Nous sommes ravis de poursuivre ce projet de collaboration internationale sans précédent et essentiel.

Le rapport jette les bases d'importantes discussions lors du Sommet d'action sur l'IA qui se tiendra en France cette année, et qui réunira des gouvernements internationaux, des entreprises de premier plan dans le domaine de l'IA, des groupes de la société civile et des experts. Ce sommet, comme le rapport, s'inscrit dans la continuité des étapes franchies lors des sommets de Bletchley Park (novembre 2023) et de Séoul (mai 2024). L'IA est l'opportunité déterminante de notre génération.

Ensemble, nous poursuivrons le dialogue et soutiendrons des actions audacieuses et ambitieuses pour maîtriser collectivement les risques de l'IA et tirer parti de ces nouvelles technologies pour le bien commun. Il n'y aura pas d'adoption de cette technologie sans sécurité : la sécurité apporte la confiance !

Nous sommes heureux de présenter ce rapport et remercions le professeur Yoshua Bengio et l'équipe de rédaction pour le travail considérable qu'ils ont consacré à son élaboration. Le Royaume-Uni et la France se réjouissent de poursuivre les discussions lors du Sommet d'action sur l'IA en février.



Clara Chappaz

Ministre délégué de la France pour
Intelligence artificielle



Le très honorable Peter Kyle, député
Secrétaire d'État britannique à la Science,
Innovation et technologie

À propos de ce rapport

- Il s'agit du premier rapport international sur la sécurité de l'IA. Après une publication intermédiaire en mai 2024, un groupe diversifié de 96 experts en intelligence artificielle (IA) ont contribué à ce premier rapport complet, dont un groupe consultatif d'experts internationaux nommés par 30 pays, l'Organisation de coopération et de développement économiques (OCDE), l'Union européenne (UE) et les Nations unies (ONU). Le rapport vise à fournir des informations scientifiques qui appuieront l'élaboration de politiques éclairées. Il ne recommande pas de politiques spécifiques.
- Le rapport est le fruit du travail d'experts indépendants. Sous la direction du Président, les experts indépendants Les rédacteurs de ce rapport ont collectivement eu toute latitude pour en décider du contenu.
- Bien que ce rapport porte sur les risques et la sécurité de l'IA, l'IA offre également de nombreux avantages potentiels Pour les particuliers, les entreprises et la société. Il existe de nombreux types d'IA, chacun présentant des avantages et des risques différents. La plupart du temps, dans la plupart des applications, l'IA aide les individus et les organisations à être plus efficaces. Mais les gens du monde entier ne pourront profiter pleinement des nombreux avantages potentiels de l'IA en toute sécurité que si ses risques sont gérés de manière appropriée. Ce rapport se concentre sur l'identification de ces risques et l'évaluation des méthodes permettant de les atténuer. Il n'a pas pour objectif d'évaluer de manière exhaustive tous les impacts sociétaux possibles de l'IA, y compris ses nombreux avantages potentiels.
- Le rapport se concentre sur l'IA à usage général. Il se concentre sur un type d'IA qui a progressé particulièrement rapidement ces dernières années et dont les risques associés ont été moins étudiés et compris : l'IA à usage général, ou l'IA capable d'effectuer une grande variété de tâches. L'analyse de ce rapport se concentre sur les systèmes d'IA à usage général les plus avancés au moment de la rédaction du présent rapport, ainsi que sur les systèmes futurs qui pourraient être encore plus performants.
- Le rapport résume les preuves scientifiques sur trois questions fondamentales : Que peut-on Que font les IA à usage général ? Quels sont les risques associés à l'IA à usage général ? Et quelles sont les techniques d'atténuation de ces risques ?
- Les enjeux sont considérables. Nous, les experts contribuant à ce rapport, continuons à être en désaccord sur plusieurs questions, mineures et majeures, concernant les capacités générales de l'IA, les risques et les mesures d'atténuation des risques. Nous considérons toutefois que ce rapport est essentiel pour améliorer notre compréhension collective de cette technologie et de ses risques potentiels. Nous espérons qu'il aidera la communauté internationale à progresser vers un plus grand consensus sur l'IA à usage général et à atténuer ses risques de manière plus efficace, afin que les gens puissent profiter en toute sécurité de ses nombreux avantages potentiels. Les enjeux sont importants. Nous sommes impatients de poursuivre cet effort.

Note du président

Point sur les dernières avancées de l'IA après la rédaction de ce rapport : note du président

Entre la fin de la période de rédaction de ce rapport (5 décembre 2024) et la publication de ce rapport en janvier 2025, un développement important a eu lieu. La société d'IA OpenAI a partagé les premiers résultats des tests d'un nouveau modèle d'IA, o3. Ces résultats indiquent des performances nettement supérieures à celles de tout modèle précédent sur un certain nombre des tests les plus difficiles du domaine de la programmation, du raisonnement abstrait et du raisonnement scientifique. Dans certains de ces tests, o3 surpasse de nombreux experts humains (mais pas tous). En outre, il réalise une percée dans un test clé de raisonnement abstrait que de nombreux experts, dont moi-même, pensaient hors de portée jusqu'à récemment. Cependant, au moment de la rédaction de cet article, il n'existe aucune information publique sur ses capacités dans le monde réel, en particulier pour résoudre des tâches plus ouvertes.

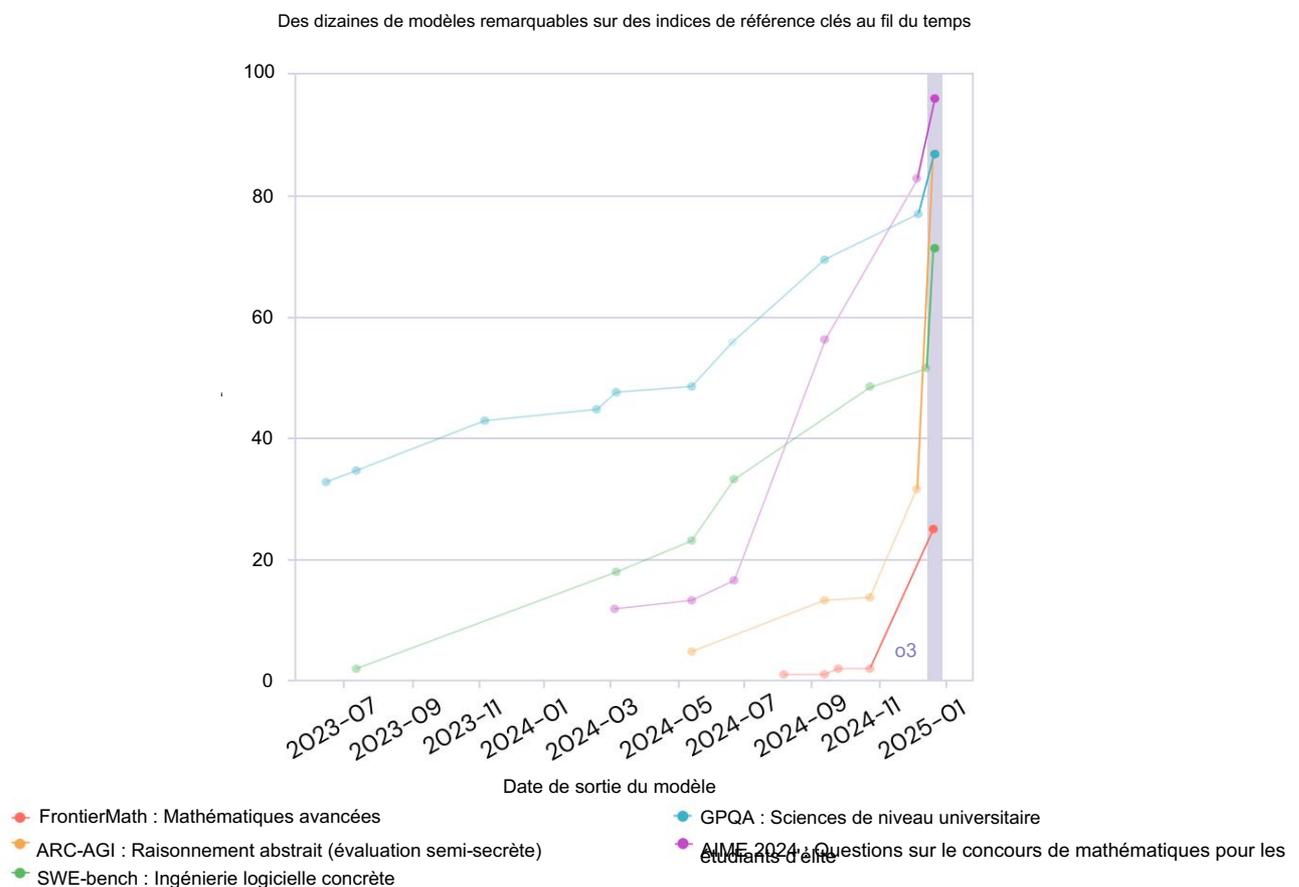


Figure 0.1 : Résultats des principaux modèles d'IA à usage général sur les principaux tests de juin 2023 à décembre 2024. o3 a montré des performances nettement améliorées par rapport à l'état de l'art précédent (zone ombrée). Ces tests de référence font partie des tests les plus difficiles du domaine de la programmation, du raisonnement abstrait et du raisonnement scientifique.

Pour le o3 non publié, la date d'annonce est indiquée ; pour les autres modèles, la date de sortie est indiquée. Certains des modèles d'IA les plus récents, dont o3, ont bénéficié d'un échafaudage amélioré et de davantage de calculs au moment des tests. Sources : Anthropic, 2024 ; Chollet, 2024 ; Chollet et al., 2025 ; Époque IA, 2024 ; Glazer et al. 2024 ; OpenAI, 2024a ; OpenAI, 2024b ;

Jimenez et al., 2024 ; Jimenez et al., 2025.

Note du président

Les résultats de l'étude o3 montrent que le rythme des progrès en matière de capacités d'IA pourrait rester élevé, voire s'accélérer. Plus précisément, ils suggèrent que donner aux modèles plus de puissance de calcul pour résoudre un problème donné (« mise à l'échelle des inférences ») peut aider à surmonter les limitations précédentes. De manière générale, la mise à l'échelle des inférences rend les modèles plus coûteux à utiliser. Mais comme l'a montré un autre modèle récent notable, R1, publié par la société DeepSeek en janvier 2025, les chercheurs travaillent avec succès à réduire ces coûts. Dans l'ensemble, la mise à l'échelle des inférences pourrait permettre aux développeurs d'IA de faire de nouvelles avancées à l'avenir. Les résultats de l'étude o3 soulignent également la nécessité de mieux comprendre comment l'utilisation croissante de l'IA par les développeurs d'IA peut affecter la vitesse du développement ultérieur de l'IA elle-même.

Les tendances mises en évidence par o3 pourraient avoir de profondes implications pour les risques liés à l'IA. Les progrès de la science et des capacités de programmation ont déjà généré davantage de preuves de risques tels que les cyberattaques et les attaques biologiques. Les résultats d'o3 sont également pertinents pour les impacts potentiels sur le marché du travail, le risque de perte de contrôle et la consommation d'énergie, entre autres. Mais les capacités d'o3 pourraient également être utilisées pour aider à se protéger contre les dysfonctionnements et les utilisations malveillantes. Dans l'ensemble, les évaluations des risques présentées dans ce rapport doivent être lues en sachant que l'IA a acquis des capacités depuis la rédaction du rapport. Cependant, jusqu'à présent, il n'existe aucune preuve des impacts réels d'o3 dans le monde réel, et aucune information ne permet de confirmer ou d'exclure des risques majeurs nouveaux et/ou immédiats.

L'amélioration des capacités suggérée par les résultats de l'étude o3 et notre compréhension limitée des implications pour les risques liés à l'IA soulignent un défi majeur pour les décideurs politiques que ce rapport identifie : ils devront souvent évaluer les avantages et les risques potentiels des avancées imminentes de l'IA sans disposer d'un large corpus de preuves scientifiques. Néanmoins, la production de preuves sur les implications en matière de sûreté et de sécurité des tendances suggérées par l'étude o3 sera une priorité urgente pour la recherche en IA dans les semaines et les mois à venir.

Principales conclusions du rapport

- Les capacités de l'IA à usage général, le type d'IA sur lequel se concentre ce rapport, ont augmenté rapidement ces dernières années et se sont encore améliorés ces derniers mois.† Il y a quelques années, les meilleurs grands modèles de langage (LLM) pouvaient rarement produire un paragraphe de texte cohérent. Aujourd'hui, l'IA à usage général peut écrire des programmes informatiques, générer des images photoréalistes personnalisées et s'engager dans des conversations ouvertes prolongées. Depuis la publication du rapport intermédiaire (mai 2024), de nouveaux modèles ont montré des performances nettement meilleures lors des tests de langage scientifique
Raisonnement et programmation.
- De nombreuses entreprises investissent désormais dans le développement d'agents d'IA à usage général, orientation potentielle pour de futurs progrès. Les agents d'IA sont des systèmes d'IA à usage général qui peuvent agir, planifier et déléguer de manière autonome pour atteindre des objectifs avec peu ou pas de surveillance humaine. Des agents d'IA sophistiqués pourraient, par exemple, utiliser des ordinateurs pour mener à bien des projets plus longs que les systèmes actuels, débloquent ainsi à la fois des avantages et des risques supplémentaires.
- De nouvelles avancées en matière de capacités dans les mois et les années à venir pourraient être lentes ou lentes. Les progrès dépendront de la capacité des entreprises à déployer rapidement davantage de données et de puissance de calcul pour former de nouveaux modèles, et de la capacité des entreprises à surmonter leurs limites actuelles en faisant évoluer les modèles de cette manière. Des recherches récentes suggèrent que l'évolution rapide des modèles pourrait rester physiquement réalisable pendant au moins plusieurs années. Mais des avancées majeures en matière de capacités peuvent également nécessiter d'autres facteurs : par exemple, de nouvelles avancées scientifiques, difficiles à prévoir, ou le succès d'une nouvelle approche d'évolution que les entreprises ont récemment adoptée.
- Plusieurs méfaits de l'IA à usage général sont déjà bien établis. Il s'agit notamment des escroqueries, Les images intimes non consensuelles (NCII) et les images d'abus sexuels sur mineurs (CSAM), les résultats des modèles biaisés à l'encontre de certains groupes de personnes ou de certaines opinions, les problèmes de fiabilité et les violations de la vie privée. Les chercheurs ont développé des techniques d'atténuation pour ces problèmes, mais jusqu'à présent, aucune combinaison de techniques ne peut les résoudre complètement. Depuis la publication du rapport intermédiaire, de nouvelles preuves de discrimination liées aux systèmes d'IA à usage général ont révélé des formes de biais plus subtiles.
- À mesure que l'IA à usage général devient plus performante, des preuves de risques supplémentaires sont progressivement mises en évidence. émergents. Il s'agit notamment de risques tels que des impacts à grande échelle sur le marché du travail, le piratage assisté par l'IA ou Les attaques biologiques et la perte de contrôle de l'IA à usage général par la société. Les experts interprètent différemment les preuves existantes sur ces risques : certains pensent que ces risques ne se manifesteront pas avant des décennies, tandis que d'autres pensent que l'IA à usage général pourrait entraîner des dommages à l'échelle de la société dans les prochaines années. Les progrès récents dans les capacités de l'IA à usage général – en particulier dans les tests de raisonnement scientifique et de programmation – ont généré de nouvelles preuves de risques potentiels tels que le piratage informatique et les attaques biologiques, ce qui a conduit une grande entreprise d'IA à revoir à la hausse son évaluation du risque biologique de son meilleur modèle de « faible » à « moyen ».

† Veuillez vous référer à la [mise à jour du Président](#) sur les dernières avancées en matière d'IA après la rédaction de ce rapport.

- Les techniques de gestion des risques sont encore balbutiantes, mais des progrès sont possibles. Il existe différentes méthodes d'évaluation et de réduction des risques liés à l'IA à usage général que les développeurs peuvent utiliser et que les régulateurs peuvent exiger, mais elles ont toutes des limites. Par exemple, les techniques d'interprétabilité actuelles permettant d'expliquer pourquoi un modèle d'IA à usage général a produit un résultat donné restent très limitées. Cependant, les chercheurs progressent dans la résolution de ces limites. En outre, les chercheurs et les décideurs politiques tentent de plus en plus de normaliser les approches de gestion des risques et de se coordonner à l'échelle internationale.
- Le rythme et l'imprévisibilité des avancées dans l'IA à usage général constituent une « preuve » Les décideurs politiques doivent souvent peser les avantages et les risques potentiels des avancées imminentes de l'IA, compte tenu des avancées parfois rapides et inattendues, sans disposer d'un large corpus de preuves scientifiques. Ce faisant, ils sont confrontés à un dilemme. D'une part, les mesures préventives d'atténuation des risques fondées sur des preuves limitées peuvent s'avérer inefficaces ou inutiles. D'autre part, attendre des preuves plus solides d'un risque imminent pourrait laisser la société démunie, voire rendre impossible l'atténuation – par exemple si des progrès soudains dans les capacités de l'IA, et les risques qui y sont associés, se produisent. Les entreprises et les gouvernements développent des systèmes d'alerte précoce et des cadres de gestion des risques qui peuvent réduire ce dilemme. Certains d'entre eux déclenchent des mesures d'atténuation spécifiques lorsqu'il existe de nouvelles preuves de risques, tandis que d'autres exigent des développeurs qu'ils fournissent des preuves de sécurité avant de lancer un nouveau modèle.
- Il existe un large consensus parmi les chercheurs sur le fait que des progrès concernant les questions suivantes Il serait utile de répondre à cette question : à quelle vitesse les capacités d'IA à usage général progresseront-elles dans les années à venir et comment les chercheurs peuvent-ils mesurer ces progrès de manière fiable ? Quels sont les seuils de risque raisonnables pour déclencher des mesures d'atténuation ? Comment les décideurs politiques peuvent-ils accéder au mieux aux informations sur l'IA à usage général qui sont pertinentes pour la sécurité publique ? Comment les chercheurs, les entreprises technologiques et les gouvernements peuvent-ils évaluer de manière fiable les risques liés au développement et au déploiement de l'IA à usage général ? Comment fonctionnent les modèles d'IA à usage général en interne ? Comment l'IA à usage général peut-elle être conçue pour se comporter de manière fiable ?
- L'IA ne nous arrive pas : les choix faits par les gens déterminent son avenir. L'avenir de L'IA à usage général est une technologie incertaine, avec un large éventail de trajectoires qui semblent possibles même dans un avenir proche, avec des résultats à la fois très positifs et très négatifs. Cette incertitude peut susciter le fatalisme et faire apparaître l'IA comme quelque chose qui nous arrive. Mais ce sont les décisions des sociétés et des gouvernements sur la façon de gérer cette incertitude qui détermineront la voie que nous emprunterons. Ce rapport vise à faciliter un débat constructif et fondé sur des données probantes sur ces décisions.

† Veuillez vous référer à la [mise à jour du Président](#) sur les dernières avancées en matière d'IA après la rédaction de ce rapport.

Résumé exécutif

Le but de ce rapport

Ce rapport synthétise l'état des connaissances scientifiques sur l'IA à usage général – une IA capable d'effectuer une grande variété de tâches – en mettant l'accent sur la compréhension et la gestion de ses risques.

Ce rapport résume les données scientifiques sur la sécurité de l'IA à usage général. L'objectif de ce rapport est de contribuer à créer une compréhension internationale commune des risques liés à l'IA avancée et de la manière dont ils peuvent être atténués. Pour y parvenir, ce rapport se concentre sur l'IA à usage général – ou l'IA capable d'effectuer une grande variété de tâches – car ce type d'IA a progressé particulièrement rapidement ces dernières années et a été largement déployé par les entreprises technologiques à diverses fins grand public et commerciales. Le rapport synthétise l'état des connaissances scientifiques sur l'IA à usage général, en mettant l'accent sur la compréhension et la gestion de ses risques.

Dans un contexte de progrès rapides, la recherche sur l'IA à usage général se trouve actuellement dans une période de découvertes scientifiques et, dans de nombreux cas, n'est pas encore une science établie. Le rapport fournit un aperçu de la compréhension scientifique actuelle de l'IA à usage général et de ses risques. Il identifie notamment les domaines de consensus scientifique et les domaines dans lesquels il existe des points de vue divergents ou des lacunes dans la compréhension scientifique actuelle.

Les citoyens du monde entier ne pourront profiter pleinement et en toute sécurité des avantages potentiels de l'IA à usage général que si les risques qu'elle comporte sont gérés de manière appropriée. Le présent rapport se concentre sur l'identification de ces risques et sur l'évaluation des méthodes techniques permettant de les évaluer et de les atténuer, notamment les moyens par lesquels l'IA à usage général peut elle-même être utilisée pour atténuer les risques. Il n'a pas pour objectif d'évaluer de manière exhaustive tous les impacts sociétaux possibles de l'IA à usage général. Plus particulièrement, les avantages actuels et futurs potentiels de l'IA à usage général – bien qu'ils soient considérables – dépassent le cadre de ce rapport. L'élaboration de politiques globales nécessite de prendre en compte à la fois les avantages potentiels de l'IA à usage général et les risques abordés dans le présent rapport. Elle nécessite également de tenir compte du fait que d'autres types d'IA présentent des profils de risques/avantages différents de ceux de l'IA à usage général actuelle.

Les trois principales sections du rapport résument les preuves scientifiques sur trois questions fondamentales :

Que peut faire l'IA à usage général ? Quels sont les risques associés à l'IA à usage général ? Et quelles sont les techniques d'atténuation de ces risques ?

Section 1 – Capacités de l'IA à usage général : Que peut faire l'IA à usage général aujourd'hui et à l'avenir ?

Les capacités de l'IA à usage général se sont rapidement améliorées ces dernières années, et les avancées futures pourraient être lentes ou extrêmement rapides.

Les capacités de l'IA contribuent largement à la plupart des risques qu'elle pose et, selon de nombreux indicateurs, les capacités de l'IA à usage général progressent rapidement. Il y a cinq ans, les principaux modèles de langage de l'IA à usage général pouvaient rarement produire un paragraphe de texte cohérent. Aujourd'hui, certains modèles d'IA à usage général peuvent engager des conversations sur un large éventail de sujets, écrire des programmes informatiques ou générer de courtes vidéos réalistes à partir d'une description. Cependant, il est techniquement difficile d'estimer et de décrire de manière fiable les capacités de l'IA à usage général.

Ces dernières années, les développeurs d'IA ont rapidement amélioré les capacités de l'IA à usage général, principalement grâce à la « mise à l'échelle ».† Ils ont continuellement augmenté les ressources utilisées pour la formation de nouveaux modèles (c'est ce que l'on appelle souvent la « mise à l'échelle ») et affiné les approches existantes pour utiliser ces ressources plus efficacement. Par exemple, selon des estimations récentes, les modèles d'IA de pointe ont connu une augmentation annuelle d'environ 4 fois des ressources de calcul (« calcul ») utilisées pour la formation et de 2,5 fois de la taille de l'ensemble de données de formation.

Le rythme des progrès futurs dans le domaine de l'IA à usage général a des implications considérables pour la gestion des risques émergents, mais les experts ne sont pas tous d'accord sur ce à quoi il faut s'attendre dans les mois et les années à venir. Les experts sont favorables à une évolution lente, rapide ou extrêmement rapide des capacités d'IA à usage général.

Les experts ne s'entendent pas sur le rythme des progrès futurs en raison de divergences de vues sur la promesse d'une « mise à l'échelle » supplémentaire. Les entreprises explorent un autre type de mise à l'échelle, nouveau et supplémentaire, qui pourrait accélérer encore les capacités.† Si la mise à l'échelle a souvent permis de surmonter les limites des systèmes précédents, les experts ne s'entendent pas sur son potentiel à résoudre les limites restantes des systèmes actuels, comme le manque de fiabilité dans l'action dans le monde physique et dans l'exécution de tâches étendues sur les ordinateurs. Ces derniers mois, un nouveau type de mise à l'échelle a montré un potentiel pour améliorer encore les capacités : plutôt que de simplement augmenter les ressources utilisées pour la formation des modèles, les entreprises d'IA s'intéressent également de plus en plus à la « mise à l'échelle par inférence », qui permet à un modèle déjà formé d'utiliser davantage de calculs pour résoudre un problème donné, par exemple pour améliorer sa propre solution, ou pour écrire ce que l'on appelle des « chaînes de pensée » qui décomposent le problème en étapes plus simples.

Plusieurs entreprises de premier plan qui développent une IA à usage général misent sur la « mise à l'échelle » pour continuer à améliorer les performances. Si les tendances récentes se poursuivent, d'ici la fin de 2026,

† Veuillez vous référer à la [mise à jour du Président](#) sur les dernières avancées en matière d'IA après la rédaction de ce rapport.

Les modèles d'IA à usage général seront formés en utilisant environ 100 fois plus de puissance de calcul que les modèles les plus gourmands en puissance de calcul de 2023, et 10 000 fois plus d'ici 2030, en combinaison avec des algorithmes qui atteignent des capacités supérieures pour une quantité donnée de calcul disponible. Outre cette mise à l'échelle potentielle des ressources de formation, les tendances récentes telles que la mise à l'échelle des inférences et l'utilisation de modèles pour générer des données de formation pourraient signifier qu'encre plus de puissance de calcul sera utilisée au total.

Il existe toutefois des obstacles potentiels à l'augmentation rapide des données et des capacités de calcul, notamment en termes de disponibilité des données, de puces d'IA, de capitaux et de capacité énergétique locale. Les entreprises qui développent une IA à usage général s'efforcent de surmonter ces obstacles potentiels.

Depuis la publication du rapport intermédiaire (mai 2024), l'IA à usage général a atteint des performances de niveau expert dans certains tests et concours de raisonnement scientifique et de programmation, et les entreprises ont déployé de gros efforts pour développer des agents d'IA autonomes.

Les progrès de la science et de la programmation ont été stimulés par des techniques de mise à l'échelle des inférences telles que l'écriture de longues « chaînes de pensée ». De nouvelles études suggèrent que la mise à l'échelle de ces approches, par exemple en permettant aux modèles d'analyser les problèmes en écrivant des chaînes de pensée encore plus longues que les modèles actuels, pourrait conduire à de nouvelles avancées dans des domaines où le raisonnement est plus important, comme la science, l'ingénierie logicielle et la planification. Outre cette tendance, les entreprises déploient de gros efforts pour développer des agents d'IA polyvalents plus avancés, capables de planifier et d'agir de manière autonome pour atteindre un objectif donné. Enfin, le prix du marché de l'utilisation d'une IA polyvalente d'un niveau de capacité donné a fortement chuté, ce qui rend cette technologie plus largement accessible et largement utilisée.

Ce rapport se concentre principalement sur les aspects techniques des progrès de l'IA, mais la vitesse à laquelle l'IA à usage général progressera n'est pas une question purement technique. Le rythme des avancées futures dépendra également de facteurs non techniques, notamment potentiellement des approches adoptées par les gouvernements pour réglementer l'IA. Ce rapport n'aborde pas la manière dont différentes approches de la réglementation pourraient affecter la vitesse de développement et d'adoption de l'IA à usage général.

Section 2 – Risques : Quels sont les risques associés à l'IA à usage général ?

Plusieurs dangers liés à l'IA à usage général sont déjà bien établis. À mesure que l'IA à usage général devient plus performante, des preuves de risques supplémentaires apparaissent progressivement.

Ce rapport classe les risques généraux liés à l'IA en trois catégories : les risques d'utilisation malveillante, les risques liés aux dysfonctionnements et les risques systémiques. Chacune de ces catégories contient des risques qui se sont déjà matérialisés ainsi que des risques qui pourraient se matérialiser dans les prochaines années.

Risques liés à une utilisation malveillante : les acteurs malveillants peuvent utiliser l'IA à usage général pour nuire à des individus, à des organisations ou à la société. Les formes d'utilisation malveillante comprennent :

- **Domage causé aux individus par le biais de faux contenus :** les acteurs malveillants peuvent actuellement utiliser l'IA à usage général pour générer des contenus factices qui nuisent aux individus de manière ciblée. Ces utilisations malveillantes comprennent la pornographie « deepfake » non consensuelle et les contenus d'abus sexuels sur mineurs générés par l'IA, la fraude financière par usurpation d'identité vocale, le chantage à des fins d'extorsion, le sabotage de réputation personnelle et professionnelle et les abus psychologiques. Cependant, bien que les rapports d'incidents faisant état de dommages causés par des contenus factices générés par l'IA soient courants, il manque encore des statistiques fiables sur la fréquence de ces incidents.
- **Manipulation de l'opinion publique :** l'IA à usage général facilite la génération d'opinions persuasives. Le contenu à grande échelle peut aider les acteurs qui cherchent à manipuler l'opinion publique, par exemple pour influencer les résultats politiques. Cependant, les preuves de la prévalence et de l'efficacité de ces efforts restent limitées. Les contre-mesures techniques telles que le tatouage numérique du contenu, bien qu'utiles, peuvent généralement être contournées par des acteurs moyennement sophistiqués.
- **Cyberattaque :** l'IA à usage général peut faciliter ou accélérer la tâche des acteurs malveillants. Les systèmes informatiques actuels ont démontré leur capacité à mener des cyberattaques de faible et moyenne complexité, et les acteurs parrainés par l'État explorent activement l'IA pour surveiller les systèmes ciblés. De nouvelles recherches ont confirmé que les capacités de l'IA à usage général liées aux cyberattaques progressent considérablement, mais il n'est pas certain que cela affectera l'équilibre entre les attaquants et les défenseurs.
- **Attaques biologiques et chimiques :** les récents systèmes d'IA à usage général ont montré une capacité à fournir des instructions et des conseils de dépannage pour reproduire des armes biologiques et chimiques connues et à faciliter la conception de nouveaux composés toxiques. Dans de nouvelles expériences visant à tester la capacité à générer des plans de production d'armes biologiques, un système d'IA polyvalent a parfois obtenu de meilleurs résultats que des experts humains ayant accès à Internet. En réponse, une société d'IA a augmenté son évaluation du risque biologique de son meilleur modèle de « faible » à « moyen ». Cependant, les tentatives concrètes de développement de telles armes nécessitent encore des ressources et une expertise supplémentaires substantielles. Une évaluation complète du risque biologique et chimique est difficile car une grande partie des recherches pertinentes sont classifiées.

Depuis la publication du rapport intermédiaire, l'IA à usage général a gagné en efficacité dans des domaines susceptibles d'être utilisés à des fins malveillantes. Par exemple, des chercheurs ont récemment mis au point des systèmes d'IA à usage général capables de détecter et d'exploiter par eux-mêmes certaines vulnérabilités en matière de cybersécurité et, avec l'aide d'un humain, de découvrir une vulnérabilité jusque-là inconnue dans des logiciels largement utilisés. Les capacités de l'IA à usage général liées au raisonnement et à l'intégration de différents types de données, qui peuvent aider à la recherche sur les agents pathogènes ou dans d'autres domaines à double usage, se sont également améliorées.

Risques liés aux dysfonctionnements : l'IA à usage général peut également causer des dommages involontaires. Même lorsque les utilisateurs n'ont pas l'intention de causer des dommages, des risques graves peuvent survenir en raison du dysfonctionnement de l'IA à usage général. Ces dysfonctionnements comprennent :

- **Problèmes de fiabilité** : l'IA à usage général actuelle peut ne pas être fiable, ce qui peut entraîner des dommages.
Par exemple, si les utilisateurs consultent un système d'IA à usage général pour obtenir des conseils médicaux ou juridiques, le système peut générer une réponse contenant des faussetés. Les utilisateurs ne sont souvent pas conscients des limites d'un produit d'IA, par exemple en raison d'une « connaissance limitée de l'IA », d'une publicité trompeuse ou d'une mauvaise communication. Il existe un certain nombre de cas connus de préjudices dus à des problèmes de fiabilité, mais les preuves sur l'ampleur exacte des différentes formes de ce problème sont encore limitées.
- **Biais** : les systèmes d'IA à usage général peuvent amplifier les biais sociaux et politiques, provoquant des
Les systèmes d'IA à usage général sont souvent biaisés en fonction de la race, du sexe, de la culture, de l'âge, du handicap, de l'opinion politique ou d'autres aspects de l'identité humaine. Cela peut conduire à des résultats discriminatoires, notamment une répartition inégale des ressources, le renforcement des stéréotypes et la négligence systématique des groupes ou des points de vue sous-représentés. Les approches techniques visant à atténuer les biais et la discrimination dans les systèmes d'IA à usage général progressent, mais elles se heurtent à des compromis entre l'atténuation des biais et des objectifs concurrents tels que l'exactitude et la confidentialité, ainsi qu'à d'autres défis.
- **Perte de contrôle** : Les scénarios de « perte de contrôle » sont des scénarios futurs hypothétiques dans lesquels une ou
Les systèmes d'IA à usage général sont de plus en plus nombreux à fonctionner en dehors du contrôle de quiconque, sans voie claire pour reprendre le contrôle. Il existe un large consensus sur le fait que l'IA à usage général actuelle n'a pas les capacités nécessaires pour poser ce risque. Cependant, l'opinion des experts sur la probabilité d'une perte de contrôle au cours des prochaines années varie considérablement : certains la considèrent comme peu plausible, d'autres la considèrent comme probable, et d'autres encore la considèrent comme un risque à probabilité modeste qui mérite d'être pris en compte en raison de sa gravité potentielle élevée. Les recherches empiriques et mathématiques en cours font progressivement avancer ces débats.

Depuis la publication du rapport intermédiaire, de nouvelles recherches ont permis d'apporter de nouvelles informations sur les risques de biais et de perte de contrôle. Les preuves de biais dans les systèmes d'IA à usage général se sont multipliées et des travaux récents ont permis de détecter d'autres formes de biais de l'IA. Les chercheurs ont observé de modestes avancées supplémentaires vers des capacités d'IA qui sont probablement nécessaires pour que les scénarios de perte de contrôle couramment évoqués se produisent. Il s'agit notamment de capacités d'utilisation autonome des ordinateurs, de programmation, d'accès non autorisé aux systèmes numériques et d'identification de moyens d'échapper à la surveillance humaine.

Risques systémiques : au-delà des risques directement liés aux capacités des modèles individuels, le déploiement généralisé de l'IA à usage général est associé à plusieurs risques systémiques plus vastes. Les exemples de risques systémiques vont des impacts potentiels sur le marché du travail aux risques pour la vie privée et aux effets environnementaux :

- **Risques liés au marché du travail** : l'IA à usage général, surtout si elle continue à progresser rapidement, a
Le potentiel d'automatisation d'un très large éventail de tâches pourrait avoir un impact significatif sur le marché du travail. Cela signifie que de nombreuses personnes pourraient perdre leur emploi actuel. Cependant, de nombreux économistes s'attendent à ce que les pertes d'emplois potentielles puissent être compensées, en partie ou potentiellement même totalement, par la création de nouveaux emplois et par une demande accrue de services non automatisés. secteurs.

- **Fracture mondiale de la R&D en IA** : la recherche et le développement (R&D) en IA à usage général sont actuellement l'IA est concentrée dans quelques pays occidentaux et en Chine. Cette « fracture de l'IA » est susceptible d'accroître la dépendance du monde à l'égard de ce petit groupe de pays. Certains experts s'attendent également à ce qu'elle contribue aux inégalités mondiales. Cette fracture a de nombreuses causes, dont certaines ne sont pas propres à l'IA. Cependant, elle résulte en grande partie des différents niveaux d'accès aux ressources informatiques très coûteuses nécessaires au développement d'une IA à usage général : la plupart des pays à revenu faible et intermédiaire (PRFI) ont un accès nettement inférieur aux ressources informatiques que les pays à revenu élevé (PRE).
- **Concentration du marché et points de défaillance uniques** : un petit nombre d'entreprises dominer le marché de l'IA à usage général. Cette concentration du marché pourrait rendre les sociétés plus vulnérables à plusieurs risques systémiques. Par exemple, si les organisations de secteurs critiques, comme la finance ou la santé, s'appuient toutes sur un petit nombre de systèmes d'IA à usage général, un bug ou une vulnérabilité dans un tel système pourrait provoquer des pannes et des perturbations simultanées à grande échelle.
- **Risques environnementaux** : utilisation croissante du calcul dans le développement de l'IA à usage général et Le déploiement de l'IA a rapidement augmenté les quantités d'énergie, d'eau et de matières premières consommées pour construire et exploiter l'infrastructure informatique nécessaire. Cette tendance ne montre aucun signe clair de ralentissement, malgré les progrès des techniques qui permettent d'utiliser l'informatique plus efficacement. L'IA à usage général a également un certain nombre d'applications qui peuvent soit bénéficier, soit nuire aux efforts de développement durable.
 - **Risques pour la vie privée** : l'IA à usage général peut provoquer ou contribuer à des violations de la vie privée des utilisateurs. Par exemple, des informations sensibles contenues dans les données d'entraînement peuvent être divulguées de manière involontaire lorsqu'un utilisateur interagit avec le système. De plus, lorsque les utilisateurs partagent des informations sensibles avec le système, ces informations peuvent également être divulguées. Mais l'IA à usage général peut également faciliter des violations délibérées de la vie privée, par exemple si des acteurs malveillants utilisent l'IA pour déduire des informations sensibles sur des individus spécifiques à partir de grandes quantités de données. Cependant, jusqu'à présent, les chercheurs n'ont pas trouvé de preuve de violations généralisées de la vie privée associées à l'IA à usage général.
- **Violations du droit d'auteur** : l'IA à usage général apprend et crée des œuvres créatives L'expression des données, remettant en cause les systèmes traditionnels de consentement, de rémunération et de contrôle des données. La collecte de données et la génération de contenu peuvent impliquer une variété de lois sur les droits des données, qui varient selon les juridictions et peuvent faire l'objet de litiges actifs. Compte tenu de l'incertitude juridique entourant les pratiques de collecte de données, les entreprises d'IA partagent moins d'informations sur les données qu'elles utilisent. Cette opacité rend la recherche sur la sécurité de l'IA tierce plus difficile.

Depuis la publication du rapport intermédiaire, de nouvelles preuves des impacts de l'IA à usage général sur le marché du travail ont émergé, tandis que de nouveaux développements ont accru les préoccupations en matière de confidentialité et de droits d'auteur. De nouvelles analyses des données du marché du travail suggèrent que les individus adoptent l'IA à usage général très rapidement par rapport aux technologies précédentes. Le rythme d'adoption par les entreprises varie considérablement selon les secteurs. En outre, les progrès récents en matière de capacités ont conduit à un déploiement croissant de l'IA à usage général dans des contextes sensibles tels que les soins de santé ou la surveillance du lieu de travail, ce qui crée de nouveaux risques pour la vie privée. Enfin, à mesure que les litiges relatifs aux droits d'auteur s'intensifient,

et les mesures techniques d'atténuation des violations du droit d'auteur restent peu fiables, les détenteurs de droits sur les données ont rapidement restreint l'accès à leurs données.

Modèles à pondération ouverte : un facteur important dans l'évaluation de nombreux risques qu'un modèle d'IA à usage général peut présenter est la manière dont il est publié au public. Les « modèles à pondération ouverte » sont des modèles d'IA dont les composants centraux, appelés « pondérations », sont partagés publiquement pour téléchargement. L'accès à une pondération ouverte facilite la recherche et l'innovation, y compris dans le domaine de la sécurité de l'IA, ainsi que l'augmentation de la transparence et la facilité pour la communauté scientifique de détecter les failles des modèles. Cependant, les modèles à pondération ouverte peuvent également présenter des risques, par exemple en facilitant une utilisation malveillante ou malavisée qu'il est difficile, voire impossible, pour le développeur du modèle de surveiller ou d'atténuer. Une fois que les pondérations du modèle sont disponibles pour téléchargement public, il n'y a aucun moyen de mettre en œuvre une restauration complète de toutes les copies existantes ou de garantir que toutes les copies existantes reçoivent des mises à jour de sécurité. Depuis le rapport intermédiaire, un consensus de haut niveau s'est dégagé sur le fait que les risques posés par une plus grande ouverture de l'IA devraient être évalués en termes de risque « marginal » : la mesure dans laquelle la publication d'un modèle à pondération ouverte augmenterait ou diminuerait un risque donné, par rapport aux risques posés par les alternatives existantes telles que les modèles fermés ou d'autres technologies.

Section 3 – Gestion des risques : Quelles techniques existe-t-il pour gérer les risques liés à l'IA à usage général ?

Plusieurs approches techniques peuvent aider à gérer les risques, mais dans de nombreux cas, les meilleures approches disponibles présentent encore des limites très importantes et ne disposent d'aucune estimation quantitative des risques ni de garanties qui sont disponibles dans d'autres domaines critiques pour la sécurité.

La gestion des risques – identifier et évaluer les risques, puis les atténuer et les surveiller – est difficile dans le contexte de l'IA à usage général. Bien que la gestion des risques ait également été très difficile dans de nombreux autres domaines, certaines caractéristiques de l'IA à usage général semblent créer des difficultés particulières.

Plusieurs caractéristiques techniques de l'IA à usage général rendent la gestion des risques dans ce domaine particulièrement difficile. Il s'agit notamment des éléments suivants :

- La gamme d'utilisations et de contextes d'utilisation possibles pour les systèmes d'IA à usage général est inhabituellement large. Par exemple, le même système peut être utilisé pour fournir des conseils médicaux, analyser le code informatique pour détecter les vulnérabilités et générer des photos. Cela augmente la difficulté d'anticiper de manière exhaustive les cas d'utilisation pertinents, d'identifier les risques ou de tester le fonctionnement des systèmes. se comportera dans des circonstances réelles pertinentes.
- Les développeurs comprennent encore peu le fonctionnement de leurs modèles d'IA à usage général. Le manque de compréhension rend plus difficile à la fois la prévision des problèmes comportementaux et l'explication et la résolution des problèmes connus une fois observés. La compréhension reste difficile principalement parce que les modèles d'IA à usage général ne sont pas programmés au sens traditionnel du terme.

Au lieu de cela, ils sont formés : les développeurs d'IA mettent en place un processus de formation qui implique un grand volume de données, et le résultat de ce processus de formation est le modèle d'IA à usage général. Le fonctionnement interne de ces modèles est en grande partie impénétrable, y compris pour les développeurs de modèles. Les techniques d'explication et d'« interprétabilité » des modèles peuvent améliorer la compréhension par les chercheurs et les développeurs du fonctionnement des modèles d'IA à usage général, mais, malgré les progrès récents, cette recherche reste naissante.

- Des agents d'IA de plus en plus performants – des systèmes d'IA à usage général qui peuvent agir de manière autonome, Les agents d'IA, qui doivent gérer les processus de manière autonome, planifier et déléguer pour atteindre les objectifs, présenteront probablement de nouveaux défis importants pour la gestion des risques. Les agents d'IA travaillent généralement à la réalisation d'objectifs de manière autonome en utilisant des logiciels généraux tels que des navigateurs Web et des outils de programmation. Actuellement, la plupart d'entre eux ne sont pas encore suffisamment fiables pour une utilisation généralisée, mais les entreprises déploient de gros efforts pour créer des agents d'IA plus performants et plus fiables et ont réalisé des progrès ces derniers mois. Les agents d'IA deviendront probablement de plus en plus utiles, mais ils pourraient également exacerber un certain nombre de risques évoqués dans ce rapport et introduire des difficultés supplémentaires pour la gestion des risques. Parmi les exemples de ces nouveaux défis potentiels, citons la possibilité que les utilisateurs ne sachent pas toujours ce que font leurs propres agents d'IA, la possibilité que les agents d'IA opèrent en dehors du contrôle de quiconque, la possibilité que des attaquants « détournent » des agents et la possibilité que les interactions entre IA créent de nouveaux risques complexes. Les approches de gestion des risques associés aux agents commencent seulement à être développées.

Outre les facteurs techniques, plusieurs facteurs économiques, politiques et autres facteurs sociétaux rendent la gestion des risques dans le domaine de l'IA à usage général particulièrement difficile.

- Le rythme des progrès de l'IA à usage général crée un « dilemme de preuves » pour Les progrès rapides des capacités permettent à certains risques d'apparaître par à-coups. Par exemple, le risque de tricherie universitaire utilisant l'IA à usage général est passé de négligeable à généralisé en l'espace d'un an. Plus un risque apparaît rapidement, plus il est difficile de le gérer de manière réactive et plus la préparation devient précieuse. Cependant, tant que les preuves d'un risque restent incomplètes, les décideurs ne peuvent pas non plus savoir avec certitude si le risque apparaîtra ou s'il est peut-être déjà apparu. Cela crée un compromis : la mise en œuvre de mesures préventives ou d'atténuation précoces peut s'avérer inutile, mais l'attente de preuves concluantes pourrait rendre la société vulnérable aux risques qui apparaissent rapidement.

Les entreprises et les gouvernements développent des systèmes d'alerte précoce et des cadres de gestion des risques qui peuvent réduire ce dilemme. Certains d'entre eux déclenchent des mesures d'atténuation spécifiques lorsqu'il existe de nouvelles preuves de risques, tandis que d'autres exigent des développeurs qu'ils fournissent des preuves de sécurité avant de commercialiser un nouveau modèle.

- Il existe un écart d'information entre ce que les entreprises d'IA savent de leurs systèmes d'IA et Ce que savent les gouvernements et les chercheurs non industriels. Les entreprises ne partagent souvent que des informations limitées sur leurs systèmes d'IA à usage général, en particulier avant leur diffusion à grande échelle. Les entreprises citent un mélange de préoccupations commerciales et de préoccupations de sécurité.

† Veuillez vous référer à la [mise à jour du Président](#) sur les dernières avancées en matière d'IA après la rédaction de ce rapport.

Il existe de nombreuses raisons de limiter le partage d'informations. Cependant, ce manque d'informations complique également la participation efficace d'autres acteurs à la gestion des risques, en particulier des risques émergents.

- Les entreprises d'IA et les gouvernements sont souvent confrontés à une forte pression concurrentielle, ce qui peut Les entreprises peuvent ainsi abandonner la gestion des risques. Dans certaines circonstances, la pression concurrentielle peut inciter les entreprises à investir moins de temps ou d'autres ressources dans la gestion des risques qu'elles ne le feraient autrement. De même, les gouvernements peuvent investir moins dans les politiques de soutien à la gestion des risques lorsqu'ils perçoivent des compromis entre la concurrence internationale et la réduction des risques.

Il existe néanmoins différentes techniques et cadres de gestion des risques liés à l'IA à usage général que les entreprises peuvent utiliser et que les régulateurs peuvent exiger. Il s'agit notamment de méthodes d'identification et d'évaluation des risques, ainsi que de méthodes d'atténuation et de surveillance de ces derniers.

- L'évaluation des risques des systèmes d'IA à usage général fait partie intégrante de la gestion des risques, mais Les évaluations de risques existantes sont très limitées. Les évaluations existantes des risques liés à l'IA à usage général reposent principalement sur des « contrôles ponctuels », c'est-à-dire sur des tests du comportement d'une IA à usage général dans un ensemble de situations spécifiques. Cela peut aider à faire apparaître les dangers potentiels avant de déployer un modèle. Cependant, les tests existants passent souvent à côté des dangers et surestiment ou sous-estiment les capacités et les risques de l'IA à usage général, car les conditions de test diffèrent du monde réel.
- Pour que l'identification et l'évaluation des risques soient efficaces, les évaluateurs ont besoin d'une expertise considérable, Les évaluateurs doivent disposer de ressources suffisantes et d'un accès suffisant aux informations pertinentes. Une évaluation rigoureuse des risques dans le contexte de l'IA à usage général nécessite de combiner plusieurs approches d'évaluation. Celles-ci vont des analyses techniques des modèles et des systèmes eux-mêmes aux évaluations des risques possibles liés à certains modes d'utilisation. Les évaluateurs doivent posséder une expertise considérable pour mener à bien ces évaluations. Pour des évaluations de risques complètes, ils ont souvent également besoin de plus de temps, d'un accès plus direct aux modèles et à leurs données de formation, et de plus d'informations sur les méthodologies techniques utilisées que celles que fournissent généralement les entreprises qui développent l'IA à usage général.
- Des progrès ont été réalisés dans la formation de modèles d'IA à usage général pour fonctionner de manière plus sûre, mais Aucune méthode actuelle ne permet d'empêcher de manière fiable des résultats même manifestement dangereux. Par exemple, une technique appelée « entraînement contradictoire » consiste à exposer délibérément des modèles d'IA à des exemples conçus pour les faire échouer ou mal se comporter pendant l'entraînement, dans le but de renforcer la résistance à de tels cas. Cependant, les adversaires peuvent toujours trouver de nouveaux moyens (« attaques ») pour contourner ces mesures de protection avec un effort faible à modéré. En outre, des preuves récentes suggèrent que les méthodes d'entraînement actuelles – qui reposent largement sur un retour d'information humain imparfait – peuvent inciter par inadvertance les modèles à induire en erreur les humains sur des questions difficiles en rendant les erreurs plus difficiles à repérer. L'amélioration de la quantité et de la qualité de ce feedback est une voie de progrès, même si les techniques de formation naissantes utilisant l'IA pour détecter les comportements trompeurs sont également prometteuses.
- Suivi – identification des risques et évaluation des performances une fois qu'un modèle est déjà utilisé – et diverses interventions visant à prévenir les actions nuisibles peuvent améliorer la sécurité d'une IA à usage général après son déploiement auprès des utilisateurs. Les outils actuels peuvent détecter les IA générées

Les systèmes d'IA peuvent détecter les contenus, suivre les performances du système et identifier les entrées/sorties potentiellement dangereuses, même si les utilisateurs moyennement qualifiés peuvent souvent contourner ces mesures de protection. Plusieurs niveaux de défense combinant des capacités de surveillance et d'intervention techniques avec une surveillance humaine améliorent la sécurité, mais peuvent entraîner des coûts et des retards. À l'avenir, les mécanismes matériels pourraient aider les clients et les régulateurs à surveiller plus efficacement les systèmes d'IA à usage général pendant le déploiement et potentiellement à vérifier les accords transfrontaliers, mais des mécanismes fiables de ce type n'existent pas encore.

- Plusieurs méthodes existent tout au long du cycle de vie de l'IA pour protéger la confidentialité. Il s'agit notamment de supprimer les méthodes de protection de la vie privée sont nombreuses et ne sont pas encore applicables aux systèmes d'IA à usage général, en raison des exigences informatiques de ces systèmes. Ces derniers mois, les méthodes de protection de la vie privée se sont développées pour répondre à l'utilisation croissante de l'IA dans des domaines sensibles, notamment les assistants pour smartphone, les agents d'IA, les assistants vocaux à l'écoute permanente et l'utilisation dans les soins de santé ou la pratique juridique.

Depuis la publication du rapport intermédiaire, les chercheurs ont fait de nouveaux progrès pour pouvoir expliquer pourquoi un modèle d'IA à usage général a produit un résultat donné.

La capacité à expliquer les décisions de l'IA pourrait contribuer à gérer les risques liés à des dysfonctionnements allant des biais et des inexactitudes factuelles à la perte de contrôle. En outre, des efforts croissants ont été déployés pour normaliser les approches d'évaluation et d'atténuation dans le monde entier.

Conclusion : Un large éventail de trajectoires pour l'avenir de l'IA à usage général est possible, et beaucoup dépendra de la manière dont les sociétés et les gouvernements agiront

L'avenir de l'IA à usage général est incertain, et de nombreuses trajectoires semblent possibles dans un avenir proche, avec des résultats à la fois très positifs et très négatifs. Mais rien n'est inéluctable dans l'avenir de l'IA à usage général. Comment l'IA à usage général sera-t-elle développée et par qui, quels problèmes elle est censée résoudre, si les sociétés seront en mesure de tirer pleinement parti du potentiel économique de l'IA à usage général, qui en bénéficie, les types de risques auxquels nous nous exposons et le montant que nous investissons dans la recherche pour gérer les risques – ces questions et bien d'autres dépendent des choix que les sociétés et les gouvernements feront aujourd'hui et à l'avenir pour façonner le développement de l'IA à usage général.

Afin de faciliter une discussion constructive sur ces décisions, ce rapport fournit un aperçu de l'état actuel de la recherche scientifique et des discussions sur la gestion des risques de l'IA à usage général.

Les enjeux sont importants. Nous sommes impatients de poursuivre cet effort.

Introduction

Nous sommes au cœur d'une révolution technologique qui va fondamentalement modifier notre façon de vivre, de travailler et d'interagir les uns avec les autres. L'intelligence artificielle (IA) promet de transformer de nombreux aspects de notre société et de notre économie.

Les capacités des systèmes d'IA se sont rapidement améliorées dans de nombreux domaines au cours des dernières années. Les grands modèles linguistiques (LLM) en sont un exemple particulièrement frappant. En 2019, GPT-2, alors le LLM le plus avancé, ne parvenait pas à produire de manière fiable un paragraphe de texte cohérent et ne pouvait pas toujours compter jusqu'à dix. Cinq ans plus tard, au moment de la rédaction de cet article, les LLM les plus puissants, tels que GPT-4, o1, Claude 3.5 Sonnet, Hunyuan-Large et Gemini 1.5 Pro, peuvent s'engager de manière cohérente dans des conversations à plusieurs tours, écrire de courts programmes informatiques, traduire entre plusieurs langues, obtenir de bons résultats aux examens d'entrée à l'université et résumer de longs documents.

Grâce à ces avancées, l'IA est désormais de plus en plus présente dans nos vies et déployée dans des contextes de plus en plus importants dans de nombreux domaines. Au cours des deux dernières années seulement, l'adoption de l'IA a connu une croissance rapide. ChatGPT, par exemple, est l'une des applications technologiques à la croissance la plus rapide de l'histoire, atteignant plus d'un million d'utilisateurs cinq jours seulement après son lancement, et 100 millions d'utilisateurs en deux mois. L'IA est désormais intégrée aux moteurs de recherche, aux bases de données juridiques, aux outils d'aide à la décision clinique et à de nombreux autres produits et services.

L'évolution des capacités et de l'adoption de l'IA, ainsi que le potentiel de progrès continu, pourraient contribuer à faire progresser l'intérêt public de plusieurs manières, mais il existe des risques. Parmi les perspectives les plus prometteuses figurent le potentiel de l'IA pour l'éducation, les applications médicales, les avancées de la recherche dans des domaines tels que la chimie, la biologie ou la physique, et l'augmentation générale de la prospérité grâce à l'innovation rendue possible par l'IA. Parallèlement à ces progrès rapides, les experts sont de plus en plus conscients des dangers actuels et des risques potentiels futurs associés aux types d'IA les plus performants.

Ce rapport vise à contribuer à une compréhension scientifique partagée au niveau international de la sécurité de l'IA avancée. Pour œuvrer à une compréhension internationale partagée des risques liés à l'IA avancée, des représentants gouvernementaux et des dirigeants du monde universitaire, des entreprises et de la société civile se sont réunis à Bletchley Park au Royaume-Uni en novembre 2023 pour le premier sommet international sur la sécurité de l'IA. Lors du Sommet, les nations présentes ont convenu de soutenir l'élaboration d'un rapport international sur la sécurité de l'IA. Ce rapport sera présenté lors du Sommet d'action sur l'IA qui se tiendra à Paris en février 2025. Une version intermédiaire de ce rapport a été publiée en mai 2024 et présentée lors du Sommet de Séoul sur l'IA. Lors du Sommet et dans les semaines et les mois qui ont suivi, les experts qui ont rédigé ce rapport ont reçu de nombreux commentaires de la part de scientifiques, d'entreprises, d'organisations de la société civile et de décideurs politiques. Ces contributions ont fortement influencé la rédaction du présent rapport, qui s'appuie sur le rapport intermédiaire et constitue le premier rapport international complet sur la sécurité de l'IA.

Un groupe international de 96 experts en IA, représentant un large éventail de points de vue et, le cas échéant, une diversité d'horizons, a contribué à ce rapport. Ils ont examiné un éventail de preuves scientifiques, techniques et socio-économiques pertinentes publiées avant le 5 décembre 2024. Le domaine de l'IA se développant rapidement, toutes les sources utilisées pour ce rapport ne sont pas évaluées par des pairs. Cependant, le rapport s'engage à ne citer que des sources de haute qualité. Les indicateurs de qualité d'une source sont les suivants :

- La pièce constitue une contribution originale qui fait avancer le domaine.
- L'article s'appuie de manière exhaustive sur la littérature scientifique existante, fait référence à la le travail des autres, le cas échéant, et l'interprète avec précision.
- L'article discute de bonne foi des objections possibles à ses revendications.
- L'article décrit clairement les méthodes employées pour son analyse. Il discute de manière critique choix des méthodes.
- L'article met clairement en évidence ses limites méthodologiques.
- La pièce a eu une influence dans la communauté scientifique.

Au moment de la rédaction du présent rapport, un consensus scientifique sur les risques liés à l'IA avancée n'était pas encore en cours d'élaboration. Dans de nombreux cas, le rapport ne présente pas d'opinions sûres. Il offre plutôt un aperçu de l'état actuel des connaissances et du consensus scientifiques, ou de leur absence. Lorsque des lacunes existent dans la littérature, le rapport les identifie, dans l'espoir que cela incitera à poursuivre les recherches.

Ce rapport ne se prononce pas sur les politiques qui pourraient être les plus appropriées pour répondre aux risques liés à l'IA. Il se veut très pertinent pour les politiques en matière d'IA, mais n'a en aucun cas valeur de prescription. En fin de compte, les décideurs politiques doivent choisir comment équilibrer les opportunités et les risques que présente l'IA avancée. Ils doivent également choisir le niveau approprié de prudence et de circonspection en réponse aux risques qui restent ambigus.

Le rapport se concentre sur l'IA « à usage général », c'est-à-dire l'IA capable d'effectuer un large éventail de tâches. L'IA est le domaine de l'informatique qui vise à créer des systèmes ou des machines capables d'effectuer des tâches qui nécessitent généralement l'intelligence humaine. Ces tâches comprennent l'apprentissage, le raisonnement, la résolution de problèmes, le traitement du langage naturel et la prise de décision. La recherche sur l'IA est un domaine d'étude vaste et en évolution rapide, et il existe de nombreux types d'IA. Ce rapport n'aborde pas tous les risques potentiels liés à tous les types d'IA avancée. Il se concentre sur l'IA à usage général, ou l'IA capable d'effectuer un large éventail de tâches. L'IA à usage général, désormais connue de beaucoup grâce à des applications telles que ChatGPT, a suscité un intérêt sans précédent pour l'IA, tant parmi le public que parmi les décideurs politiques, au cours des deux dernières années. Les capacités de l'IA à usage général se sont améliorées particulièrement rapidement. L'IA à usage général est différente de ce qu'on appelle « l'IA étroite », un type d'IA spécialisé pour effectuer une tâche spécifique ou quelques tâches très similaires.

Pour mieux comprendre comment ce rapport définit l'IA à usage général, il est utile de faire une distinction entre les « modèles d'IA » et les « systèmes d'IA ». Les modèles d'IA peuvent être considérés comme l'essence mathématique brute qui est souvent le « moteur » des applications d'IA. Un système d'IA est une combinaison de plusieurs

composants, y compris un ou plusieurs modèles d'IA, conçus pour être particulièrement utiles aux humains d'une manière ou d'une autre. Par exemple, l'application ChatGPT est un système d'IA ; son moteur principal, GPT-4, est un modèle d'IA.

Le rapport couvre les risques liés à la fois aux modèles d'IA à usage général et aux systèmes d'IA à usage général. Aux fins du présent rapport :

- Un modèle d'IA est un modèle d'IA à usage général s'il peut exécuter, ou peut être adapté pour exécuter, une grande variété de tâches. Si un tel modèle est adapté pour effectuer principalement un ensemble plus restreint de tâches, il est toujours considéré comme un modèle d'IA à usage général.
- Un système d'IA est un système d'IA à usage général s'il est basé sur un modèle d'IA à usage général.

« Adapter un modèle » fait ici référence à l'utilisation de techniques telles que le réglage fin d'un modèle (l'entraînement d'un modèle déjà pré-entraîné sur un ensemble de données nettement plus petit que l'ensemble de données précédent utilisé pour l'entraînement), son incitation de manière spécifique (« ingénierie rapide ») et des techniques d'intégration du modèle dans un système plus large.

Les grands modèles et systèmes d'IA génératifs, tels que les chatbots basés sur des LLM, sont des exemples bien connus d'IA à usage général. Ils permettent une génération flexible de résultats qui peuvent facilement s'adapter à un large éventail de tâches distinctes. L'IA à usage général comprend également les IA qui peuvent effectuer un large éventail de tâches suffisamment distinctes dans un domaine spécifique tel que la biologie structurale.

Dans le domaine de l'IA à usage général, ce rapport se concentre sur l'IA à usage général qui est au moins aussi performante que l'IA à usage général la plus avancée d'aujourd'hui. Parmi les exemples, citons GPT-4o, AlphaFold-3 et Gemini 1.5 Pro. Notez que dans la définition de ce rapport, un modèle ou un système n'a pas besoin d'avoir plusieurs modalités (par exemple, la parole, le texte et les images) pour être considéré comme à usage général.

Ce qui compte, c'est la capacité à réaliser une grande variété de tâches, ce qui peut également être accompli par un modèle ou un système avec une seule modalité.

L'IA à usage général ne doit pas être confondue avec l'« intelligence artificielle générale » (IAG). Le terme IAG n'a pas de définition universelle, mais il est généralement utilisé pour désigner une future IA potentielle qui égale ou dépasse les performances humaines dans toutes ou presque toutes les tâches cognitives. En revanche, plusieurs modèles et systèmes d'IA actuels répondent déjà aux critères permettant de les considérer comme de l'IA à usage général, tels que définis dans ce rapport.

Ce rapport ne traite pas des risques liés à l'IA « étroite », qui est formée pour effectuer une tâche spécifique et capture un corpus de connaissances très limité. L'accent mis sur l'IA polyvalente avancée est dû au fait que les progrès dans ce domaine ont été les plus rapides et que les risques associés sont moins étudiés et compris. L'IA étroite peut toutefois également être très pertinente du point de vue des risques et de la sécurité, et les preuves relatives aux risques de ces systèmes sont utilisées dans tout le rapport.

Les modèles et systèmes d'IA étroits sont utilisés dans une vaste gamme de produits et de services dans des domaines tels que la médecine, la publicité ou la banque, et peuvent présenter des risques importants. Ces risques peuvent entraîner des préjudices tels que des décisions d'embauche biaisées, des accidents de voiture ou des recommandations de traitement médical néfastes. IA étroite

Le terme « armes létales autonomes » est également utilisé dans diverses applications militaires, par exemple dans les systèmes d'armes létales autonomes (SALA) (1). Ces sujets sont traités dans d'autres forums et ne relèvent pas du champ d'application de ce rapport. La portée des futurs rapports potentiels n'est pas encore déterminée.

Un groupe important et diversifié d'experts internationaux de premier plan a contribué à ce rapport, notamment des représentants désignés par 30 pays de tous les groupes régionaux des Nations Unies, ainsi que par l'OCDE, l'UE et l'ONU. Bien que nos points de vue diffèrent parfois, nous partageons la conviction qu'un débat scientifique et public constructif sur l'IA est nécessaire pour que les populations du monde entier puissent profiter en toute sécurité des avantages de cette technologie. Nous espérons que ce rapport pourra contribuer à ce débat et servir de base à de futurs rapports qui amélioreront progressivement notre compréhension commune des capacités et des risques de l'IA avancée.

Le rapport est organisé en cinq sections principales : Après cette introduction, 1. Capacités de l'IA à usage général fournit des informations sur les capacités actuelles de l'IA à usage général, les principes sous-jacents et les tendances futures potentielles. 2. Risques examine les risques associés à l'IA à usage général. 3. Approches techniques de la gestion des risques présente les approches techniques permettant d'atténuer les risques liés à l'IA à usage général et évalue leurs points forts et leurs limites. La conclusion résume et conclut.

1. Capacités de IA à usage général

1.1. Comment l'IA à usage général est-elle développée

INFORMATIONS CLÉS

- L'IA à usage général peut effectuer et aider les utilisateurs à accomplir une grande variété de tâches.
par exemple, il peut produire du texte, des images, des vidéos, de l'audio, des actions ou des annotations pour des données.
- L'IA à usage général est basée sur l'apprentissage profond. L'apprentissage profond exploite de grandes quantités de ressources informatiques permettant à un modèle d'IA d'apprendre des modèles utiles à partir d'une grande quantité de données de formation.
- Le cycle de vie d'une IA à usage général peut être divisé en étapes distinctes. Ces étapes sont les suivantes :
 - Collecte et prétraitement des données : les développeurs et les travailleurs des données collectent, nettoient, étiquettent, normaliser et transformer les données de formation brutes en un format à partir duquel le modèle peut apprendre efficacement.
 - Pré-formation : les développeurs alimentent les modèles d'IA avec de vastes quantités de données pour inculquer des connaissances générales. Connaissance par l'apprentissage à partir d'exemples. Cette étape est actuellement celle qui nécessite le plus de calculs.
 - Ajustement : les développeurs et les travailleurs des données sous contrat affinent davantage les données pré-formées « modèle de base » dans un processus appelé « réglage fin » pour optimiser les performances du modèle pour une application spécifique ou le rendre plus utile en général. Cette étape peut demander beaucoup de travail.
 - Intégration système : les développeurs combinent un ou plusieurs modèles d'IA à usage général avec d'autres composants, tels que des interfaces utilisateur ou des filtres de contenu, pour améliorer les capacités et la sécurité et pour produire un « système d'IA » complet prêt à l'emploi.
 - Déploiement : les développeurs mettent le système d'IA intégré à la disposition d'autres personnes pour qu'elles l'utilisent. mettre en œuvre le système d'IA dans des applications ou des services du monde réel.
 - Suivi post-déploiement : les développeurs recueillent et analysent les commentaires des utilisateurs, suivent des mesures de performance et apportent des améliorations itératives pour résoudre les problèmes ou les limitations découverts lors d'une utilisation réelle. Ces améliorations peuvent inclure des ajustements plus précis ou une mise à jour de l'intégration du système.
- Depuis la publication du rapport intermédiaire (mai 2024), les capacités de l'IA à usage général
Les tests de raisonnement en plusieurs étapes se sont améliorés. Cela est dû en grande partie à des techniques de réglage fin grâce auxquelles un modèle apprend à aborder les problèmes de manière plus structurée avant de générer un résultat.

Définitions clés

- **Modèle** : Un programme informatique, souvent basé sur l'apprentissage automatique, conçu pour traiter les entrées et générer des résultats. Les modèles d'IA peuvent effectuer des tâches telles que la prédiction, la classification, la prise de décision ou la génération, constituant ainsi le cœur des applications d'IA.

- **Système** : une configuration intégrée qui combine un ou plusieurs modèles d'IA avec d'autres composants, comme des interfaces utilisateur ou des filtres de contenu, pour produire une application avec laquelle les utilisateurs peuvent interagir avec.
- **Calcul** : abréviation de « ressources informatiques », qui fait référence au matériel (par exemple unités de traitement graphique (GPU), logiciels (par exemple, logiciels de gestion de données) et infrastructures (par exemple, centres de données) nécessaires pour former et exécuter les systèmes d'IA.
- **Apprentissage profond** : une technique d'apprentissage automatique dans laquelle de grandes quantités de données et de calculs sont utilisés pour former des réseaux neuronaux artificiels multicouches (inspirés du cerveau biologique) afin d'apprendre et d'extraire automatiquement des fonctionnalités de haut niveau à partir de grands ensembles de données, permettant ainsi de puissantes capacités de reconnaissance de modèles et de prise de décision.
- **Développeur** : Toute organisation qui conçoit, construit, intègre, adapte ou combine des modèles d'IA ou des systèmes.
- **Réseau neuronal** : Un type de modèle d'IA constitué d'une structure mathématique inspirée
Le cerveau humain est composé de nœuds interconnectés (comme des neurones) qui traitent et apprennent à partir des données. Les systèmes d'IA à usage général actuels sont basés sur des réseaux neuronaux.
- **Pondérations** : paramètres du modèle qui représentent la force de la connexion entre les nœuds d'un réseau neuronal. Les poids jouent un rôle important dans la détermination de la sortie d'un modèle en réponse à une entrée donnée et sont mis à jour de manière itérative pendant la formation du modèle pour améliorer son performance.

L'« IA à usage général » fait référence aux modèles ou systèmes d'intelligence artificielle qui peuvent effectuer un large éventail de tâches plutôt que d'être spécialisés dans une fonction spécifique. Alors que toute IA fonctionne sur une base fondamentale d'entrée-sortie (traitement des données pour générer des résultats), l'IA à usage général se distingue par sa capacité à gérer un large éventail de tâches, par exemple résumer du texte, générer des images ou écrire du code informatique (pour une définition plus détaillée de l'IA à usage général, voir [Introduction](#)). Cette polyvalence la rend utile, permettant des applications dans de nombreux domaines tels que [la santé](#), [la finance](#) et [l'ingénierie](#). Cependant, ces capacités présentent également de nouveaux défis, notamment en matière de sécurité et d'utilisation éthique. La complexité de la gestion de plusieurs cas d'utilisation potentiels augmente le risque de conséquences imprévues, de biais et d'utilisation abusive.

Voici quelques exemples d'IA à usage général :

- Modèles de langage, tels que o1 (2*), GPT-4o (3*), Gemini-1.5 (4*), Claude-3.5 (5*), Command r+ (6*), Qwen2.5 (7*), la famille ERNIE (8*), Hunyuan-Large (9*), Yi-Lightning (10*), Llama-3.1 (11*) et Mistral Large (12*).
- Générateurs d'images (13), tels que DALL-E 3 (14*) et Stable Diffusion-3 (15*).
- Générateurs de vidéos tels que SORA (16*), Pika (17) et Runway (17).
- Systèmes de robotique et de navigation, tels que PaLM-E (18) et Octo (19*).
- Des agents d'IA capables d'accomplir des tâches relativement complexes dans la poursuite d'un objectif avec peu d'effort humain implication, comme AutoGPT (20), Sibyl (21*) et « The AI Scientist » (22*).
- Prédicteurs de structures biomoléculaires, tels que AlphaFold-3 (23).

Les modèles d'IA à usage général sont développés via un processus appelé « apprentissage profond ». L'apprentissage profond est un paradigme de développement de l'IA axé sur la création de systèmes informatiques qui apprennent à partir d'exemples. Au lieu de programmer des règles spécifiques dans les systèmes, les chercheurs alimentent ces systèmes avec des exemples – Les individus apprennent à reconnaître des modèles et à donner un sens à de nouvelles informations, comme des images, des textes ou des sons. L'apprentissage profond a commencé à émerger comme paradigme dominant du développement de l'IA au début des années 2010. Il s'est consolidé comme paradigme principal après des développements notables tels que la victoire du système AlphaGo contre le meilleur joueur mondial de Go en 2016.

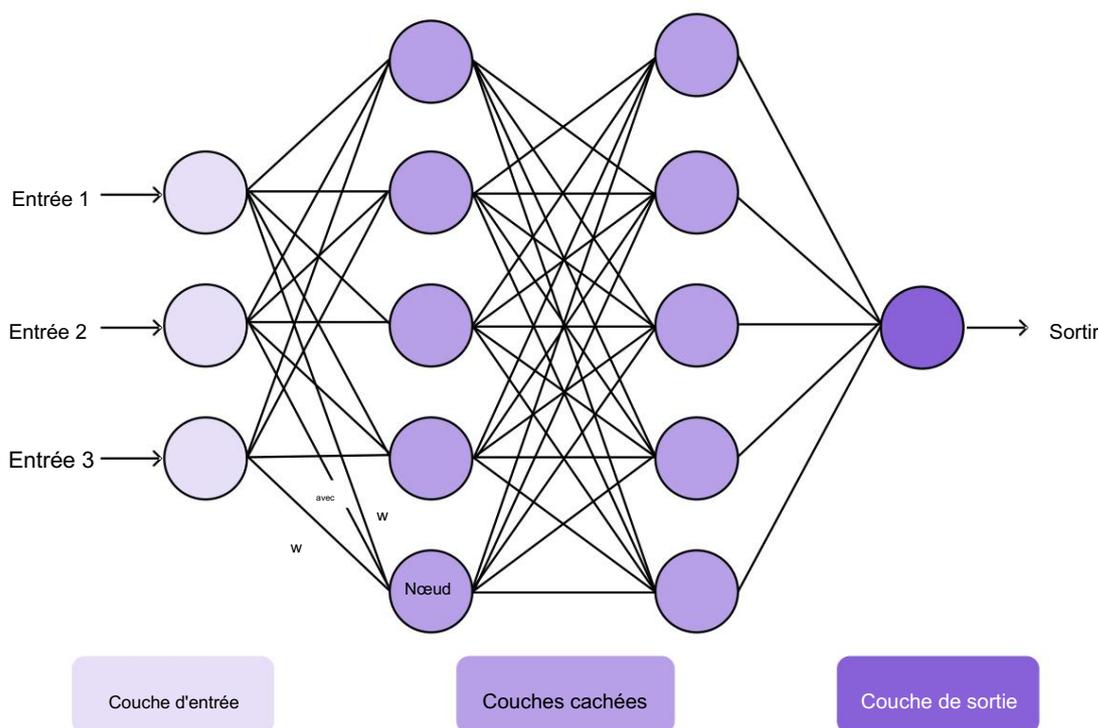


Figure 1.1 : Les modèles d'IA à usage général actuels sont des réseaux neuronaux, inspirés du cerveau animal. Ces réseaux sont composés de nœuds connectés, où la force des connexions entre les nœuds est appelée « poids ». Les poids sont mis à jour par un entraînement itératif avec de grandes quantités de données. Source : International AI Safety Report.

Il existe de nombreux types d'IA polyvalente, mais ils sont développés à l'aide de méthodes et de principes communs. L'apprentissage profond fonctionne en traitant les données via des « couches » de nœuds mathématiques interconnectés (voir la figure 1.1), souvent appelés « neurones » car ils sont vaguement inspirés des neurones du cerveau biologique (« réseaux neuronaux ») (24). À mesure que l'information circule d'une couche de neurones à la suivante, le modèle affine ses représentations. Par exemple, dans un système de vision, les premières couches peuvent détecter des caractéristiques simples telles que des contours ou des formes de base dans une image, tandis que les couches plus profondes combinent ces caractéristiques pour reconnaître des motifs plus complexes comme des visages ou des objets. Lorsque le système fait des erreurs, les algorithmes d'apprentissage profond ajustent la force des différentes connexions entre les neurones pour améliorer les performances du modèle. La force de chaque connexion entre les nœuds est souvent appelée « poids ». Cette approche de l'apprentissage en couches est à l'origine du nom de l'apprentissage profond, et elle est efficace pour des tâches qui nécessitaient auparavant une intelligence humaine. La plupart des modèles d'IA à usage général de pointe reposent désormais sur une architecture de réseau neuronal spécifique connue sous le nom de « Transformer » (25), capable de traiter simultanément de grandes quantités de données.

ont été très efficaces pour apprendre à partir de grandes quantités de données, ce qui a conduit à des améliorations significatives dans la traduction et la génération de textes et a finalement conduit au développement de LLM tels que GPT-4o.

Le processus de développement et de déploiement de l'IA à usage général suit une série d'étapes distinctes.

Ces étapes se déroulent à des moments différents, dépendent de ressources différentes, requièrent des techniques différentes et sont parfois réalisées par des développeurs différents (voir la figure 1.2 / le tableau 1.1). Par conséquent, les différentes politiques et réglementations affectant les données, les ressources informatiques (« calcul ») ou la surveillance humaine peuvent affecter chaque étape différemment.

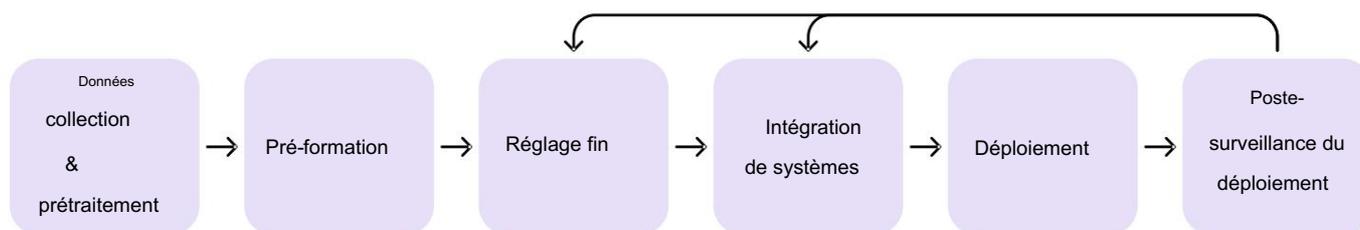


Figure 1.2 : Le processus de développement et de déploiement d'une IA à usage général suit une série d'étapes distinctes, depuis la collecte et le prétraitement des données jusqu'à la surveillance post-déploiement. Source : International AI Safety Report.

Avant de former un modèle d'IA à usage général, les développeurs collectent et préparent des données appropriées, ce qui constitue une opération à grande échelle. La création d'ensembles de données de formation de haute qualité implique des pipelines complexes de collecte, de nettoyage et de conservation des données. Les ensembles de données de formation derrière les modèles de pointe comprennent un nombre immense d'exemples provenant de partout sur Internet. Les équipes développent souvent des systèmes de filtrage sophistiqués pour réduire le contenu inapproprié ou nuisible, éliminer les données en double et améliorer la représentation sur différents sujets et perspectives. Le prétraitement des données peut également aider à réduire les problèmes de droits d'auteur et de confidentialité, à gérer plusieurs langues et formats et à améliorer la documentation sur la provenance des données. De nombreuses entreprises emploient de grandes équipes d'annotateurs et d'experts en la matière pour vérifier et étiqueter des parties des données, développer des systèmes de classification pour la qualité du contenu et créer des ensembles de données spécialisés pour des capacités spécifiques.

Collecte de données et pré-traitement	Les développeurs collectent, nettoient, étiquettent, normalisent et transforment les données d'entraînement brutes en un format à partir duquel le modèle peut apprendre. Il s'agit d'un processus très exigeant en main-d'œuvre.
Pré-formation	Les développeurs alimentent les modèles avec des quantités massives de données diverses (texte, code, images, etc.) pour leur inculquer des connaissances générales. La préformation produit un « modèle de base ». Il s'agit d'un processus très gourmand en ressources de calcul.
Réglage fin	Les développeurs continuent à former le modèle de base pour l'optimiser pour une application spécifique ou le rendre plus utile en général. Cela se fait généralement à l'aide d'une grande quantité de commentaires générés par l'homme. Il s'agit d'un processus moyennement intensif en calcul et très exigeant en main-d'œuvre.

Intégration système	Les développeurs combinent un ou plusieurs modèles d'IA à usage général avec d'autres composants tels que des interfaces utilisateur ou des filtres de contenu pour créer un « système d'IA » complet prêt à l'emploi.
Déploiement	Les développeurs mettent le système d'IA intégré à la disposition d'autres personnes.
Surveillance post-déploiement	Les développeurs collectent et analysent les commentaires des utilisateurs, suivent les mesures d'impact et de performance et apportent des améliorations itératives pour résoudre les problèmes ou les limitations découverts lors d'une utilisation réelle.

Tableau 1.1 : À chaque étape du cycle de vie de l'IA, le modèle d'IA est amélioré pour une utilisation en aval et finalement déployé en tant que système d'IA entièrement intégré.

Lors de la pré-formation, les développeurs présentent des modèles d'IA à usage général avec de grandes quantités de données, ce qui permet au modèle d'apprendre des modèles. Au début du processus de formation, un modèle non formé produit des résultats aléatoires. Cependant, grâce à l'exposition à des millions ou des milliards d'exemples – Des images, des textes ou des fichiers audio permettent au modèle d'apprendre progressivement des faits et des modèles qui lui permettent de donner un sens aux informations dans leur contexte. La préformation produit un « modèle de base » doté de connaissances et de capacités générales.

La préformation des modèles d'IA à usage général est souvent l'étape de développement la plus intensive en termes de calcul. Le processus de préformation prend des semaines ou des mois et utilise des dizaines de milliers d'unités de traitement graphique (GPU) ou d'unités de traitement tensoriel (TPU), des puces informatiques spécialisées conçues pour traiter rapidement de nombreux calculs. Aujourd'hui, ce processus utilise environ 10 milliards de fois plus de calcul que la formation de modèles de pointe en 2010 (26). Certains développeurs effectuent la préformation avec leur propre calcul, tandis que d'autres utilisent des ressources fournies par des fournisseurs de calcul spécialisés. Dans tous les cas, les coûts énergétiques sont élevés et il est prévu que pour les plus grands modèles d'IA à usage général, les coûts de calcul de préformation à eux seuls dépasseront 1 milliard de dollars pour certains modèles d'ici 2027 (27). Voir [2.3.4. Risques pour l'environnement](#) pour une discussion sur les coûts environnementaux de la formation.

Après une pré-formation, les modèles d'IA à usage général apprennent à partir de commentaires spécialement organisés et d'ensembles de données spécialisés pour améliorer les performances et l'efficacité du modèle – un processus appelé « réglage fin ». Après la préformation, la plupart des modèles d'IA à usage général subissent une ou plusieurs étapes de réglage supplémentaires pour affiner leur capacité à accomplir les tâches prévues. Le réglage fin peut inclure diverses techniques, notamment l'apprentissage à partir d'exemples souhaitables (28, 29) ou de renforcement positif/négatif (30, 31*). D'une certaine manière, le réglage fin d'une IA à usage général peut être comparé à l'enseignement d'un élève par la pratique et le retour d'information. Souvent, le réglage fin suit ce schéma : 1. les chercheurs donnent à un modèle de base des tâches qu'il essaie ensuite de résoudre ; 2. les chercheurs marquent ensuite les bonnes réponses comme des exemples positifs et les erreurs comme des exemples négatifs ; 3. le modèle est ensuite mis à jour de telle sorte qu'il tende à privilégier les approches qui ont bien fonctionné et à éviter celles qui n'ont pas fonctionné, devenant progressivement plus fiable. Dans l'ensemble, le réglage fin améliore les performances des modèles d'IA à usage général en leur permettant d'utiliser les connaissances et les capacités existantes pour accomplir la tâche souhaitée. Le réglage fin est traditionnellement l'étape de formation la plus exigeante en main-d'œuvre,

Les systèmes d'IA à usage général sont de plus en plus utilisés pour affiner d'autres modèles à usage général (32*, 33*). En pratique, l'ajustement est généralement un processus itératif dans lequel les développeurs alternent entre l'ajustement et les tests jusqu'à ce que leurs tests montrent que le système répond aux spécifications souhaitées.

Après le réglage fin vient l'« intégration système », qui consiste à combiner des modèles d'IA à usage général avec d'autres composants tels que des interfaces utilisateur ou des filtres de contenu pour produire un système d'IA à usage général. Un système d'IA à usage général est une combinaison d'un ou plusieurs modèles d'IA à usage général et de tous les composants supplémentaires nécessaires pour les rendre opérationnels, tels que des interfaces utilisateur, une infrastructure de traitement des données et divers outils. Par exemple, GPT-4o est un modèle d'IA à usage général qui traite du texte, des images et de l'audio. Cependant, ChatGPT est un système d'IA à usage général qui combine le modèle GPT-4o avec une interface de chat, un traitement de contenu, un accès Web et une intégration d'applications pour créer un produit fonctionnel. Les composants supplémentaires d'un système d'IA visent également à améliorer les capacités, l'utilité et la sécurité. Par exemple, un système peut être équipé d'un filtre qui détecte et bloque les entrées ou les sorties de modèle contenant du contenu nuisible. Les développeurs conçoivent également de plus en plus ce que l'on appelle des « échafaudages » autour de modèles d'IA à usage général qui leur permettent de planifier à l'avance, de poursuivre des objectifs et d'interagir avec le monde (voir [1.2. Capacités actuelles](#)). Tout comme le réglage fin, l'intégration du système implique généralement une alternance d'étapes d'intégration et de test. La dernière étape avant le déploiement consiste généralement à établir un rapport sur le développement, les capacités et les résultats des tests du système. C'est ce que l'on appelle souvent une « carte système » (34).

Après l'intégration du système, le « déploiement » permet de mettre à disposition les systèmes d'IA. Le déploiement est le processus d'implémentation des systèmes d'IA dans des applications, des produits ou des services du monde réel où ils peuvent répondre aux demandes et fonctionner dans un contexte plus large. Le déploiement peut prendre plusieurs formes : déploiement interne pour une utilisation par le développeur du système, ou déploiement externe, soit publiquement, soit pour des clients privés. On sait très peu de choses publiquement sur les déploiements internes. Cependant, on sait que les entreprises adoptent différents types de stratégies pour le déploiement externe. Par exemple, les entreprises offrent souvent un accès via des interfaces utilisateur en ligne ou des intégrations qui permettent à leurs modèles d'être utilisés avec des applications personnalisées conçues par des développeurs en aval. Ces intégrations peuvent permettre aux systèmes d'IA à usage général d'un développeur d'être utilisés dans de nombreuses autres applications. Par exemple, une entreprise peut concevoir un chatbot de service client sur mesure alimenté par le système d'IA à usage général d'une autre entreprise.

Le « déploiement » et la « diffusion de modèles » sont des activités distinctes qui sont facilement confondues. Le « déploiement » implique la mise en service d'un système d'IA intégré comme décrit ci-dessus. La « diffusion de modèles », en revanche, implique la mise à disposition de modèles formés pour que les entités en aval puissent les utiliser, les étudier, les modifier et/ou les intégrer dans leurs propres systèmes. Il existe un éventail d'options de diffusion de modèles allant de la fermeture complète à l'ouverture complète (35*). Les modèles entièrement fermés sont conservés uniquement pour la recherche et le développement internes. Les modèles entièrement ouverts sont ceux pour lesquels tous les composants du modèle (par exemple, les poids, le code, les données de formation) et la documentation sont mis à disposition gratuitement sous une licence open source pour que chacun puisse les utiliser, les étudier, les partager ou les modifier (36*). Certaines IA à usage général

Certains modèles, comme GPT-4o (3*), se situent à l'extrémité fermée du spectre, tandis que d'autres se situent plutôt à l'extrémité ouverte du spectre. Par exemple, Llama-3.1 (37*) a des pondérations « ouvertes » qui sont disponibles en téléchargement public. Du point de vue de l'atténuation des risques, les formes plus ouvertes de diffusion des modèles présentent des avantages et des inconvénients (voir [2.4. Impact des modèles d'IA polyvalents à pondération ouverte sur les risques liés à l'IA](#)).

Après le déploiement, les développeurs peuvent se lancer dans une « surveillance » – inspecter les entrées et les sorties du système pour suivre les performances et détecter les problèmes – et mettre à jour leurs systèmes en continu. Ce processus implique la collecte et l'analyse des commentaires des utilisateurs, le suivi des mesures de performance et la réalisation d'améliorations itératives pour résoudre les problèmes ou les limitations découverts lors d'une utilisation dans le monde réel (38). Ces améliorations peuvent inclure un réglage plus précis ou une mise à jour de l'intégration du système. Dans la pratique, il existe souvent un « jeu du chat et de la souris » dans lequel les développeurs mettent continuellement à jour les systèmes de haut niveau en réponse aux problèmes nouvellement découverts (39). Voir [3.4.2. Surveillance et intervention](#), pour une discussion sur les méthodes de surveillance des systèmes d'IA à usage général et d'intervention si nécessaire.

Depuis la publication du rapport intermédiaire, les développeurs ont réalisé des progrès significatifs dans les techniques d'intégration de systèmes qui peuvent permettre à l'IA à usage général d'effectuer un raisonnement plus avancé. En septembre 2024, OpenAI a annoncé son nouveau modèle prototype o1 avec des méthodes d'échafaudage et de formation plus avancées qui ont permis des gains de performance significatifs sur des tâches telles que les mathématiques et la programmation (2*). Contrairement aux modèles précédents, o1 utilise la résolution de problèmes par « chaîne de pensée » qui décompose les problèmes en étapes qui sont ensuite résolues petit à petit. La chaîne de pensée a permis des améliorations dans les tâches complexes - o1 a obtenu 83 % aux examens de qualification de l'Olympiade internationale de mathématiques (IMO) contre 13 % pour GPT-4o - et est considérée comme une étape importante vers le développement d'agents d'IA : des systèmes d'IA à usage général qui peuvent interagir de manière autonome avec le monde, planifier à l'avance et poursuivre des objectifs. Cependant, le processus amélioré de résolution de problèmes nécessite beaucoup plus de temps et de calcul à la fois pendant la formation et au point d'utilisation. L'étendue du raisonnement et les capacités du modèle restent floues (40).

Français Les décideurs politiques sont confrontés à divers défis découlant de la manière dont l'IA à usage général est développée. Des risques et des vulnérabilités peuvent apparaître à de nombreux moments du processus de développement et de déploiement, ce qui rend les interventions les plus efficaces difficiles à identifier et à hiérarchiser. Les progrès dans le développement de modèles se produisent également rapidement et sont difficiles à prévoir. Il est donc difficile d'articuler des interventions politiques solides qui vieilliront bien avec une technologie en évolution rapide. Non seulement les risques et les vulnérabilités associés à l'IA à usage général sont susceptibles de changer, mais les exigences de développement de modèles le sont également. Par exemple, les modèles basés sur le raisonnement tels que o1 nécessitent des ressources de calcul beaucoup plus importantes au point d'utilisation, ce qui présente de nouvelles implications pour la planification de l'infrastructure de calcul à long terme. [1.2. Capacités actuelles](#) et [1.3. Capacités dans les années à venir](#) élargissent l'état des capacités actuelles de l'IA et la manière dont ces capacités sont susceptibles d'évoluer, ce qui pose de nouveaux risques et défis.

1.2. Capacités actuelles

INFORMATIONS CLÉS†

- Comprendre et mesurer les capacités de l'IA à usage général est essentiel pour évaluer leurs risques. Les cadres de gouvernance et les engagements existants reposent sur une mesure précise des capacités d'IA à usage général, mais elles constituent une cible mouvante et difficile à mesurer et à définir.
- La plupart des experts s'accordent à dire que les systèmes d'IA à usage général sont capables d'effectuer des tâches telles que :
 - Assister les programmeurs et réaliser des projets d'ingénierie logicielle de petite et moyenne envergure tâches.
 - Créer des images difficiles à distinguer des vraies photographies.
 - Participer à une conversation fluide dans de nombreuses langues.
 - Trouver et résumer des informations pertinentes à une question ou à un problème à partir de nombreuses données sources.
 - Travailler simultanément avec plusieurs « modalités » telles que le texte, la vidéo et la parole.
 - Résoudre des problèmes de mathématiques et de sciences de manuels scolaires jusqu'au niveau universitaire.
- La plupart des experts s'accordent à dire que l'IA à usage général n'est actuellement pas capable d'effectuer des tâches telles que :
 - Effectuer des tâches robotiques utiles telles que des travaux ménagers.
 - Éviter systématiquement les fausses déclarations.
 - Exécuter de manière autonome des projets de longue durée, tels que des programmes ou des recherches sur plusieurs jours projets.
- Les agents d'IA à usage général peuvent de plus en plus agir et planifier de manière autonome en contrôlant Les principales entreprises d'IA investissent massivement dans les agents d'IA, car elles estiment qu'ils sont rentables. Les tests liés à la navigation sur le Web, au codage et aux tâches de recherche progressent rapidement, même si les agents d'IA actuels ont encore du mal à effectuer des tâches qui nécessitent de nombreuses étapes.
- Depuis la publication du rapport intermédiaire (mai 2024), les systèmes d'IA polyvalents ont considérablement amélioré leurs performances dans les tests de raisonnement scientifique et de programmation. Ces améliorations proviennent en partie de techniques qui permettent à l'IA polyvalente de décomposer des problèmes complexes en étapes plus petites, en écrivant ce que l'on appelle des « chaînes de pensée », avant de les résoudre.
- L'un des principaux défis pour les décideurs politiques est de savoir comment tenir compte des capacités spécifiques au contexte dans les réglementations. Les capacités de l'IA à usage général peuvent changer considérablement avec des réglages plus précis, des incitations et des outils mis à disposition du système. Elles peuvent également diminuer dans des contextes inconnus. Des évaluations plus rigoureuses sont nécessaires pour éviter de surestimer ou de sous-estimer les capacités.

† Veuillez vous référer à la [mise à jour du Président](#) sur les dernières avancées en matière d'IA après la rédaction de ce rapport.

Définitions clés

- **Modalités** : les types de données qu'un système d'IA peut recevoir avec compétence en entrée et produire en sortie, notamment du texte (langage ou code), des images, des vidéos et des actions robotiques.
- **Capacités** : l'éventail des tâches ou des fonctions qu'un système d'IA peut exécuter et comment il peut les exécuter avec compétence.
- **Améliorations du temps d'inférence** : techniques utilisées pour améliorer les performances d'un système d'IA après son entraînement initial, sans modifier le modèle sous-jacent. Cela comprend des méthodes d'invite intelligentes, des méthodes de sélection de réponses (par exemple, l'échantillonnage de plusieurs réponses et le choix d'une réponse majoritaire), l'écriture de longues « chaînes de pensée », l'« échafaudage » des agents, etc.
- **Échafaudage** : logiciel supplémentaire construit autour d'un système d'IA qui l'aide à effectuer une tâche. Par exemple, un système d'IA peut avoir accès à une application de calcul externe pour améliorer ses performances sur des problèmes arithmétiques. Un échafaudage plus sophistiqué peut structurer les résultats d'un modèle et guider le modèle pour améliorer ses réponses étape par étape.
- **Chaîne de pensée** : Un processus de raisonnement dans lequel une IA génère des étapes intermédiaires ou explications lors de la résolution d'un problème ou de la réponse à une question. Cette approche imite le raisonnement logique humain et la délibération interne, aidant le modèle à décomposer les tâches complexes en étapes séquentielles plus petites pour améliorer la précision et la transparence de ses résultats.
- **Inférence** : Le processus dans lequel une IA génère des sorties basées sur une entrée donnée, appliquer les connaissances acquises lors de la formation.
- **Agent IA** : une IA à usage général qui peut élaborer des plans pour atteindre des objectifs, effectuer de manière adaptative des tâches impliquant plusieurs étapes et des résultats incertains en cours de route, et interagir avec son environnement (par exemple en créant des fichiers, en effectuant des actions sur le Web ou en déléguant des tâches à d'autres agents) avec peu ou pas de surveillance humaine.
- **Évaluations** : Évaluations systématiques des performances, des capacités et des performances d'un système d'IA. vulnérabilités ou impacts potentiels. Les évaluations peuvent inclure des analyses comparatives, des équipes rouges et des audits et peuvent être menées avant et après le déploiement du modèle.
- **Benchmark** : un test ou une mesure standardisé, souvent quantitatif, utilisé pour évaluer et comparer les performances des systèmes d'IA sur un ensemble fixe de tâches conçues pour représenter le monde réel usage.

Cette section se concentre sur les capacités de base des modèles et systèmes d'IA à usage général qui sont aujourd'hui disponibles au public. La section [1.3. Capacités dans les années à venir, traite des développements](#) futurs attendus dans les capacités de l'IA, et la section [2. Risques](#), traite des capacités dangereuses spécifiques et de leurs applications associées qui contribuent aux risques.

Les capacités d'un système d'IA à usage général sont difficiles à mesurer de manière fiable (41). Une mise en garde importante concernant les évaluations des capacités de l'IA est que leurs profils de capacités et la cohérence avec laquelle ils présentent certaines capacités diffèrent considérablement de ceux des humains. Par exemple, deux études révèlent que les modèles de langage échouent plus souvent sur les problèmes de comptage et d'arithmétique impliquant des nombres qui sont rares dans leurs données d'entraînement (42*, 43). Le succès d'un système d'IA à un test de capacités dépend fortement des exemples particuliers choisis pour le test, ainsi que de la façon dont il est

demandé ou chargé de les résoudre (ce qui, en pratique, dépend de la compétence de son utilisateur) – ce qui rend particulièrement difficile de garantir l'absence d'une capacité dans un système d'IA (par exemple une capacité qui pourrait entraîner des risques sociétaux (44*)) ; voir [2.1 Risques liés à une utilisation malveillante](#). La diversité des données et un investissement approprié dans des méthodes permettant d'obtenir le comportement souhaité d'un modèle (par exemple par le biais d'améliorations du temps d'inférence telles que l'échafaudage, l'incitation et le réglage fin) peuvent contribuer à rendre l'évaluation des capacités plus fiable.

Modalités d'entrée et de sortie

Les « modalités » d'un système d'IA sont les types de données qu'il peut utilement recevoir en entrée et produire en sortie. Par exemple, les systèmes d'IA polyvalents avec une modalité texte peuvent prendre en compte le texte saisi par l'utilisateur ou des documents sources, et produire un langage naturel cohérent, engager des conversations et répondre à des questions de compréhension de lecture sur un passage. Les systèmes d'IA avec des modalités image et texte peuvent être capables de répondre à des questions sur le contenu des images ou de générer des images selon des instructions en langage naturel. Comprendre les modalités qu'un système d'IA polyvalent peut traiter est important pour développer une intuition sur les grands ensembles de tâches qu'il pourrait être capable d'accomplir en théorie, et les menaces possibles que lui-même – et les futurs modèles de ce type – peuvent représenter. Il existe des systèmes polyvalents pour plus de 9 modalités (45)

Les systèmes d'IA polyvalents sont de plus en plus utilisés dans le traitement de textes, d'audio, d'images et de vidéos, certains systèmes se concentrant spécifiquement sur une modalité supplémentaire telle que les actions robotiques, les représentations de protéines et d'autres molécules, les données de séries chronologiques (46*) ou la musique (47*). Cependant, les systèmes de traitement de texte et d'images tels que ChatGPT sont à l'origine d'une grande partie de l'attention actuelle portée à l'IA polyvalente. Les systèmes d'IA polyvalents avancés sont de plus en plus capables de traiter des entrées et de générer des sorties dans de multiples modalités telles que le texte, la vidéo et la parole.

Texte et code : les systèmes d'IA à usage général peuvent dialoguer de manière interactive et écrire de courts programmes informatiques. Les modèles linguistiques avancés peuvent générer du texte et dialoguer de manière interactive dans une variété de langues naturelles, de sujets et de formats. On peut citer comme exemples GPT-4 d'OpenAI, Claude d'Anthropic et Gemini de Google, ainsi que des modèles disponibles en libre accès de Meta (la série de modèles Llama), Mistral AI, Alibaba (la série Qwen) et DeepSeek (48*, 49*, 50*, 51*, 52*, 53*, 54*). En plus du langage humain, ces modèles peuvent traiter et générer de nombreux types de données codées sous forme de texte, notamment des formules mathématiques et du code informatique. Ils peuvent écrire des programmes de courte à moyenne durée, aider les développeurs de logiciels et effectuer des actions informatiques (telles que des recherches sur le Web) lorsqu'ils disposent de moyens tels qu'un accès à Internet (55, 56).

Audio et parole : les systèmes d'IA à usage général peuvent engager une conversation orale et imiter de manière convaincante la voix humaine. Certains systèmes d'IA à usage général, notamment GPT-4o (3*) et Gemini 1.5 (49*), peuvent traiter l'audio de la même manière que le texte, en répondant à des questions sur le contenu d'un clip audio (par exemple, une conversation parlée). Une étude récente sur l'utilisation de l'IA étroite pour la synthèse de texte à la parole a révélé que sur deux tests de synthèse vocale universitaires, la voix d'une personne pouvait être reproduite de manière convaincante dans un son de haute qualité à partir d'une durée de seulement trois secondes

enregistrement (57*). Le système d'IA à usage général GPT-4o peut converser en temps réel avec un discours de type humain dans son « mode vocal avancé » et peut émuler une variété de voix humaines.

Images : Les systèmes d'IA à usage général peuvent décrire le contenu des images avec une grande précision, générer des images selon une description détaillée et effectuer d'autres tâches basées sur des images. De nombreux systèmes d'IA à usage général peuvent utiliser des images à la fois comme entrée et comme sortie. Les systèmes d'IA à usage général tels que Claude, et GPT-4o, Pixtral et Qwen2-VL peuvent décrire le contenu des images en langage, y compris les objets et les activités qui y sont représentés (3*, 50*, 58*, 59*, 60*). Les modèles les plus performants peuvent donner un sens à des images et des documents complexes, Anthropic signalant que son système Claude 3.5 Sonnet peut répondre correctement à plus de 90 % des questions dans trois tests de référence qui impliquent le traitement de documents, de graphiques et de diagrammes scientifiques, représentant des paramètres de tests standardisés humains (5*). Les systèmes d'IA à usage général peuvent également générer des images en sortie, avec un contenu et un style spécifiés en langage humain (par exemple, des systèmes tels que Stable Diffusion 3 (15*) et DALL-E 3 (14*)). Les progrès réalisés dans les modèles de génération d'images permettent de contrôler plus facilement le contenu et le style des images, de représenter des scènes de plus en plus complexes et réalistes et de produire des images quasiment impossibles à distinguer des images naturelles (14*). D'autres systèmes d'IA polyvalents peuvent effectuer des tâches basées sur les images, telles que la catégorisation des objets représentés dans les images (61) et l'identification de leur emplacement (62*).

Vidéo : Les systèmes d'IA à usage général peuvent transcrire ou décrire le contenu de vidéos et générer de courtes vidéos selon des instructions, mais le mouvement représenté dans ces vidéos n'est pas toujours réaliste. Certains systèmes d'IA à usage général peuvent prendre une vidéo en entrée et analyser son contenu, comme V-JEPA (63*), Gemini 1.5 (49*), GPT-4o (3*) et Qwen2-VL (60*). Ces systèmes peuvent permettre de rechercher et d'analyser un contenu long, par exemple en localisant des moments clés ou des éléments d'information révélés dans une vidéo. Certains systèmes à usage général peuvent également générer des vidéos réalistes en haute définition, par exemple Sora (16*) et Movie Gen (64*). Ces modèles peuvent générer de courtes vidéos (moins d'une minute) représentant une scène décrite dans du texte, éventuellement en référence à d'autres images et vidéos. Ils peuvent modifier des vidéos selon des instructions (par exemple en changeant la saison représentée de l'été à l'hiver) et générer des vidéos représentant des individus sur des photographies de référence (par exemple en effectuant une activité décrite). Ces vidéos semblent généralement réalistes, même si la précision des scènes générées par rapport au texte d'instruction a tendance à être moins bonne que pour les systèmes de génération d'images de pointe, et les vidéos contiennent souvent des mouvements non naturels ou physiquement impossibles qui les distinguent clairement des vidéos naturelles. Les modèles vidéo avancés ne sont arrivés sur le marché qu'en 2024 et leurs implications sont encore à l'étude.

Actions robotiques : les systèmes d'IA à usage général peuvent être utilisés pour planifier des mouvements robotiques, mais ne peuvent pas encore contrôler eux-mêmes des robots ou des machines physiques. Les systèmes d'IA à usage général peuvent être utilisés pour planifier des actions robotiques en plusieurs étapes et traduire le langage pédagogique en action robotique plans (65*, 66). Les chercheurs explorent également des modèles d'IA à usage général qui non seulement planifient ou interprètent, mais génèrent également des actions robotiques, comme le RT-2-X de Google (67), et la société de conduite autonome Waymo développe des modèles d'IA à usage général pour générer des plans de conduite et des modèles de l'environnement d'un véhicule (68*). Cependant, les capacités des modèles d'IA à usage général à générer

Les actions robotiques sont relativement rudimentaires. Cela s'explique en partie par le fait que la collecte de données pour les actions nécessite généralement l'utilisation de robots physiques et qu'elle est difficile à réaliser à très grande échelle (69), bien que des efforts substantiels soient déployés (67, 70, 71). Les systèmes d'IA à usage général ne peuvent pas encore contrôler efficacement les robots ou les machines physiques pour effectuer de nombreuses tâches utiles telles que les travaux ménagers, car l'intégration de modèles d'IA à usage général avec des systèmes de contrôle moteur reste un défi (72).

Protéines et autres molécules : les systèmes d'IA à usage général peuvent effectuer une gamme de tâches utiles aux biologistes, telles que la prédiction du repliement des protéines et l'aide à la conception des protéines. Les systèmes d'IA à usage général qui fonctionnent avec des protéines et d'autres grosses molécules fonctionnent à l'aide de diverses représentations (par exemple, séquences de résidus, structures 3D). Ces modèles peuvent prédire les structures protéiques dans diverses conditions (par exemple, dans des complexes protéine-protéine), générer de nouvelles protéines utiles et effectuer un large éventail de tâches liées aux protéines pertinentes pour la découverte et la conception de médicaments (73*), les qualifiant de « modèles de base » (74) et de modèles d'IA à usage général selon la définition de ce rapport (voir [Introduction](#)). Ils peuvent de plus en plus être utilisés pour générer des conceptions de nouvelles protéines avec des fonctions prévisibles dans de grandes familles de protéines (75, 76).

Améliorations après la pré-formation

Les tâches qu'un système d'IA à usage général peut accomplir dépendent des techniques qui lui sont appliquées après la pré-formation initiale. Une étude de 16 méthodes d'amélioration a révélé qu'elles nécessitent généralement moins de 1 % des ressources informatiques nécessaires à la mise en œuvre de la pré-formation des systèmes, tout en améliorant les capacités de ces systèmes à peu près autant que ce que l'on pourrait attendre en consacrant 5 fois plus de ressources à la pré-formation (77). Cela suggère que la politique de développement et de déploiement de systèmes d'IA à usage général peut devoir anticiper l'effet que ces améliorations auront sur les capacités des systèmes d'IA à usage général. Certaines méthodes d'amélioration courantes (77, 78) comprennent :

- Réglage fin : Le réglage fin fait référence à la formation supplémentaire du modèle de base pré-entraîné pour l'optimiser pour une application spécifique ou la rendre plus utile de manière générale, par exemple en l'entraînant à suivre des instructions.
- Améliorations du temps d'inférence : l'inférence est le processus par lequel un modèle d'IA génère des sorties basées sur une entrée donnée, appliquant ainsi les connaissances acquises pendant la formation. Les améliorations du temps d'inférence sont une classe de techniques d'intégration de systèmes qui modifient les entrées d'un modèle et organisent ses sorties. Les exemples incluent la production de plusieurs réponses candidates à une question et la sélection de la meilleure d'entre elles (79*, 80*), la production de longues « chaînes de pensée » (voir le paragraphe suivant) pour résoudre des problèmes complexes (2*), ou l'utilisation d'hybrides de ces approches (81). D'autres améliorations du temps d'inférence incluent :
 - Méthodes d'incitation : élaborer les instructions du système pour améliorer ses performances, par exemple en lui fournissant des exemples de problèmes et de solutions (82, 83), en lui fournissant des documents utiles pour le contexte ou en lui demandant de « penser étape par étape » (84) ;
 - Échafaudage d'agent et utilisation d'outils : fournir au modèle les moyens de décomposer une tâche de haut niveau en un plan avec des sous-objectifs clairs et de déléguer à des copies de lui-même

exécuter chaque étape du plan en interagissant avec son environnement, par exemple en utilisant des sites Web (85*) ou exécuter du code (86*, 87, 88*) pour effectuer son travail en tant qu'agent IA (89, 90).

Les modèles d'IA à usage général se sont nettement améliorés pour répondre aux questions scientifiques de niveau doctorat

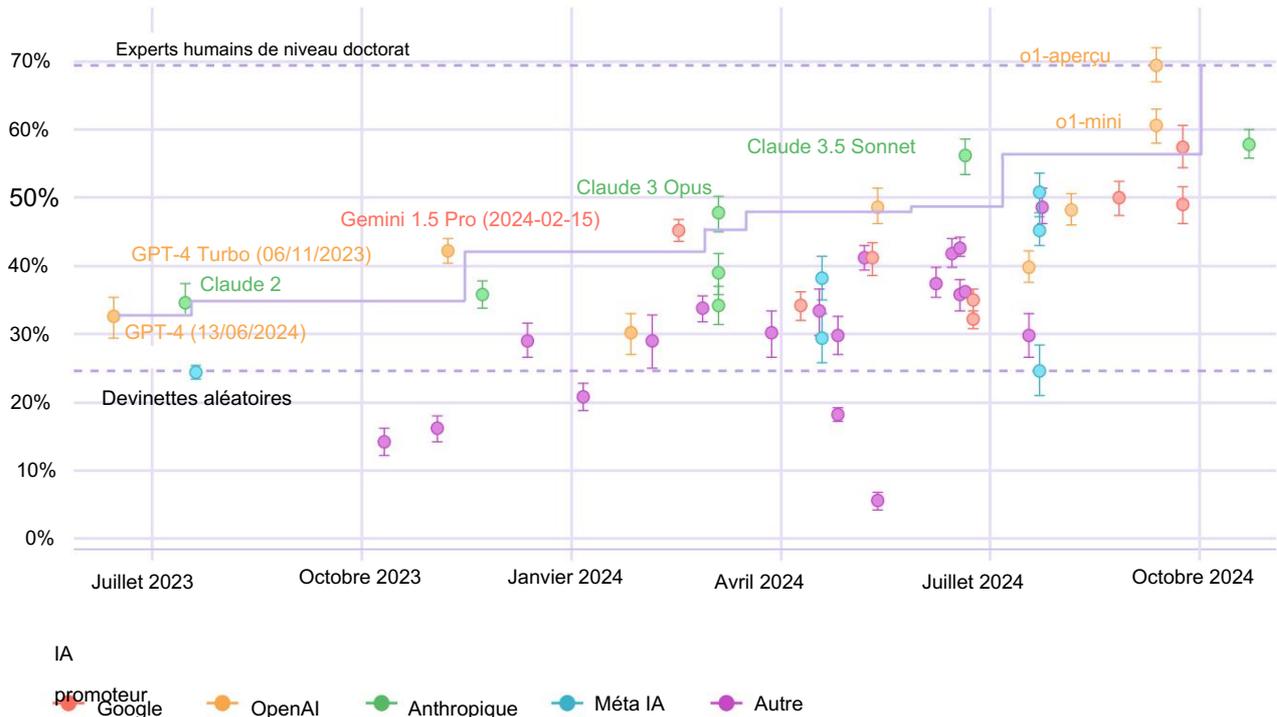


Figure 1.3 : Depuis la publication du rapport intermédiaire (mai 2024), les modèles d'IA à usage général ont connu une augmentation rapide de leurs performances pour répondre aux questions scientifiques de niveau doctorat. Les chercheurs ont testé des modèles sur GPQA Diamond, un ensemble de questions à choix multiples difficiles sur la biologie, la chimie et la physique, auxquelles les personnes sans expertise de niveau doctorat dans chaque domaine ne sont pas en mesure de répondre correctement même avec un accès à Internet. Sur ces tests, la précision est passée de 33 % avec GPT-4 en juin 2023 (légèrement au-dessus des devinettes aléatoires) à 49 % avec GPT-4o en mai 2024, atteignant 70 % (en faisant correspondre des experts titulaires d'un doctorat dans le domaine de chaque question) avec o1-preview en septembre 2024. Cette augmentation est en partie due au fait qu'o1-preview écrit une longue « chaîne de pensée » dans laquelle il peut décomposer le problème et essayer différentes approches avant de produire sa réponse finale. Pour les progrès réalisés sur d'autres tests, voir la figure 1.4 dans [1.3 Capacités dans les années à venir](#). Source : Epoch AI, 2024 (91).

Depuis la publication du rapport intermédiaire, des études ont montré que les capacités d'un système d'IA à usage général peuvent être considérablement augmentées en lui faisant consacrer plus de temps et de calculs à chaque problème individuel. Le système o1 d'OpenAI, lancé en septembre 2024, a obtenu un score suffisamment élevé à l'American Invitational Mathematics Examination (AIME) pour se qualifier pour l'Olympiade mathématique des États-Unis, et a atteint un niveau de doctorat expert sur des questions de physique, de chimie et de biologie de niveau postuniversitaire organisées pour une difficulté élevée (92*)

(voir Figure 1.3). La clé des améliorations d'o1 a été de tirer parti de calculs supplémentaires au moment de l'inférence en écrivant une longue « chaîne de pensée » pour décomposer le problème et travailler sur des hypothèses. Une autre amélioration populaire du temps d'inférence tire parti d'un calcul accru pendant le temps d'inférence en échantillonnant plusieurs sorties du modèle et en choisissant parmi elles. Deux études récentes menées par des chercheurs de l'industrie, du milieu universitaire et de la société civile examinent comment

Les capacités augmentent avec la quantité de calcul en temps d'inférence utilisant de telles techniques (93, 94*). Ils ont constaté que les capacités augmentent à un rythme qui est approximativement logarithmique avec l'investissement en temps de calcul d'inférence, formant une tendance similaire à la relation entre la croissance des capacités et le temps de calcul de formation comme décrit dans [la section 1.3. Capacités dans les années à venir](#). Ceci, combiné au succès [d'GPT-1](#), suggère que la quantité de calcul en temps d'inférence consacrée à chaque problème pourrait être un levier à usage général par lequel les capacités d'un système d'IA à usage général existant peuvent être augmentées (en particulier dans les applications scientifiques et technologiques), c'est-à-dire en lui permettant simplement de produire une « chaîne de pensée » beaucoup plus longue avant sa réponse. Cependant, susciter des capacités améliorées en utilisant plus de calcul en temps d'inférence nécessite plus de calcul, ce qui augmente les coûts.

Que peut faire l'IA polyvalente actuelle ?

Les modèles de langage à usage général peuvent répondre correctement à de nombreuses questions de bon sens et factuelles, mais ils peuvent être incohérents et commettre des erreurs insignifiantes. Les systèmes d'IA à usage général codent une large gamme de faits, les systèmes de pointe actuels obtenant en moyenne un score supérieur à 92 % aux tests de connaissances de niveau licence dans des matières telles que la chimie et le droit (92*).

Cependant, ces systèmes ne parviennent souvent pas à identifier des distinctions factuelles subtiles ou des arguments contradictoires (95, 96), sont susceptibles de fournir des réponses biaisées sur la base des modèles d'interaction des utilisateurs (97, 98), sont moins précis pour répondre aux questions sur des scénarios inhabituels (42*, 99*, 100) et génèrent généralement des citations, biographies ou faits totalement inexistantes ou faux (101, 102*, 103, 104, 105), ou font de simples erreurs de bon sens (106, 107). Ces problèmes sont interprétés par certains chercheurs comme indiquant qu'ils manquent d'une véritable compréhension du fonctionnement du monde (108) et rendent difficile l'adoption de tels systèmes dans des environnements qui nécessitent une grande fiabilité. Voir [1.3. Capacités dans les années à venir pour une discussion plus approfondie](#).

Les systèmes d'IA à usage général peuvent atteindre des performances similaires ou supérieures à celles des experts humains sur certaines tâches autonomes de connaissances et de raisonnement, mais ils commettent toujours des erreurs sur des problèmes simples d'une manière dont les humains ne le font pas. Dans une étude, un système d'IA à usage général a pu prédire la probabilité d'événements futurs avec une précision rivalisant avec celle des prévisionnistes experts sur une plateforme de prévision en ligne (109). En ce qui concerne le codage, GPT-1 se situe au 89e percentile des humains sur Codeforces, une plateforme de codage compétitive en ligne, et peut résoudre 41 % d'un échantillon de tâches d'ingénierie autonomes et réelles tirées de la plateforme de partage de code GitHub (2*).

Cependant, même sur des problèmes mathématiques simples de niveau primaire, les systèmes d'IA à usage général présentent des schémas d'erreur différents de ceux des humains. Par exemple, deux études montrent que leur précision diminue considérablement lorsque des phrases manifestement non pertinentes sont insérées dans le problème. (110*, 111*), avec une réduction de 17,5 % de la précision pour une version préliminaire de GPT-1 (110*). Deux études récentes révèlent également que lorsque les systèmes d'IA à usage général sont confrontés à des problèmes qui nécessitent davantage d'étapes de raisonnement pour être résolus, leur taux d'erreur augmente plus rapidement que ce à quoi on pourrait s'attendre s'ils avaient un taux d'erreur constant par étape (110*, 112*). Cela suggère que les systèmes d'IA à usage général ne peuvent pas être utilisés pour des problèmes complexes, et conduit certains chercheurs à affirmer que ces systèmes « ne peuvent pas effectuer un véritable raisonnement logique » (110*), bien que les avis sur ce point parmi les experts soient mitigés.

Des études montrent que l'assistance de l'IA rend les développeurs de logiciels plus productifs, et l'adoption d'outils d'IA pour la programmation est en hausse. Des études sur GitHub Copilot, une aide populaire au codage de l'IA, montrent des gains de productivité allant de 8-22 % (113) à 56 % (114*). Les développeurs se perçoivent comme étant plus productifs lorsqu'ils sont interrogés (115), et l'assistance de l'IA est généralement plus bénéfique pour les développeurs inexpérimentés (114*, 115). Dans une enquête menée auprès de plus de 65 000 développeurs de logiciels de mai à juin 2024 par Stack Overflow, un forum communautaire de questions-réponses sur la programmation populaire, 63 % des développeurs de logiciels professionnels ont déclaré utiliser des outils d'IA dans leur flux de travail (116), contre 44 % l'année précédente (117).

Depuis la publication du rapport intermédiaire, les agents d'IA à usage général qui exécutent indépendamment des tâches sur l'ordinateur ont fait l'objet d'investissements importants et deviennent rapidement plus fiables sur les critères de référence conçus pour tester le potentiel d'automatisation du travail. Les agents d'IA sont des systèmes d'IA à usage général qui peuvent de manière autonome élaborer des plans, exécuter des tâches complexes et interagir avec leur environnement en contrôlant des logiciels et des ordinateurs, avec peu de surveillance humaine. Les agents d'IA peuvent être créés en équipant les systèmes d'IA à usage général d'une fine couche de logiciel supplémentaire appelée « échafaudage ». Les exemples de tâches pour les agents d'IA comprennent les tâches de navigation sur le Web telles que répondre à des questions (85*) ou faire des achats en ligne (118, 119), l'assistance à la recherche scientifique (22*, 120, 121*), le développement de logiciels (122), la formation de modèles d'apprentissage automatique (123*, 124, 125*, 126), la réalisation de cyberattaques (127), le suivi d'instructions pour naviguer dans des environnements simulés (128) ou le contrôle de robots physiques (19*). Français Sur la plupart de ces tâches, les agents d'IA actuels réussissent dans les cas de complexité faible à moyenne, mais échouent lorsque la tâche nécessite de nombreuses étapes ou devient plus complexe. Dans une étude d'évaluation portant sur 77 tâches, allant de tâches simples telles que l'exploitation des vulnérabilités de base des sites Web à des tâches complexes en plusieurs étapes telles que la formation de modèles d'apprentissage automatique, des modèles de pointe tels que GPT-4o, o1 et Claude 3.5 Sonnet ont réussi près de 40 % des tâches lorsqu'ils étaient équipés d'un échafaudage d'agent, un taux similaire à celui des humains qui sont limités à 30 minutes pour chaque tâche (2*, 129). Dans la même étude, o1 a fait quelques progrès - sans réussir complètement - sur deux des sept tâches difficiles conçues pour refléter les tâches difficiles de la recherche et du développement (R&D) en IA, telles que l'optimisation du code du réseau neuronal (2*, 129). Les progrès dans ce domaine sont rapides : de nouvelles architectures d'agents sont rapidement développées (130*, 131*, 132*), et le taux de réussite du système le plus performant sur un sous-ensemble de haute qualité de SWE-bench, un benchmark d'agent d'ingénierie logicielle populaire, est passé de 22 % à 45 % d'avril à août 2024 (122).

Depuis la publication du rapport intermédiaire, les chercheurs ont également progressé dans l'exploitation de nouveaux types de données multimodales pour former des modèles d'IA pour le contrôle des robots. Une approche consiste à former un système sur un grand ensemble de données de vidéos annotées avec des descriptions textuelles de leur contenu, suivi d'un ensemble de données plus petit de vidéos (rares) annotées avec des commandes d'action du robot (133*). Une deuxième nouvelle approche utilise l'IA polyvalente activée par la vision existante pour traduire des vidéos d'humains en plans d'action pour les robots, et forme des modèles de contrôle des robots à l'aide de ces données (134). Une troisième nouvelle approche s'entraîne uniquement sur la vidéo, mais implique que les modèles apprennent implicitement les actions qui y sont décrites, ce qui permet au modèle de s'adapter rapidement au contrôle de nouveaux robots, même si sa formation initiale ne portait que sur des vidéos d'humains (135*). Ces études suggèrent

que de nouvelles méthodes exploitant l'apprentissage multimodal ouvriront bientôt le goulot d'étranglement des données qui empêche actuellement les développeurs de former des systèmes d'IA à usage général pour contrôler les robots.

Les principales lacunes en matière de données probantes concernant les capacités actuelles de l'IA sont les suivantes :

- Il n'existe pas d'index complet et constamment mis à jour des capacités de l'IA.

Les capacités de l'IA deviennent rapidement obsolètes à mesure que de nouveaux modèles sont publiés et que des améliorations du temps d'inférence sont développées. La compréhension des capacités de l'IA par les chercheurs progresse grâce à un patchwork relativement ponctuel de publications universitaires et industrielles qu'il peut être difficile de synthétiser pour obtenir une image complète. Idéalement, les décideurs politiques devraient avoir accès à des preuves actualisées, fiables, standardisées et complètes.

- Les évaluations des capacités de l'IA ne sont souvent pas répliquées sur de nouvelles données. Études d'évaluation fournir des exemples d'un système d'IA effectuant une tâche (ou ne parvenant pas à le faire) sur certains échantillons de données, mais ces exemples ne se reproduisent souvent pas lorsque les expériences sont réexécutées ou essayées sur des données différentes (136). Pour que les évaluations soient fiables et reproductibles, elles doivent idéalement être exécutées sur des ensembles de données vastes et diversifiés, qui sont étendus au fil du temps.

- Il n'existe pas de normes communes pour mesurer la manière dont l'IA augmente les capacités humaines.

Il n'existe pas encore de critères de référence normalisés pour « l'amélioration » – mesurant l'efficacité avec laquelle les humains peuvent utiliser des systèmes d'IA à usage général pour accomplir diverses tâches, par rapport à l'utilisation de la technologie existante – qui peuvent informer le public de cet aspect du progrès.

Des tests sont effectués – bien que les détails soient souvent confidentiels – pour les risques d'utilisation abusive de substances chimiques, biologiques, radiologiques et nucléaires (CBRN) ; voir [2.1.4. Attaques biologiques et chimiques](#) et [2.4 Impact des modèles d'IA polyvalents à pondération ouverte sur les risques d'IA.](#)

Pour les décideurs politiques, les principaux défis sont les suivants :

- Les mesures standardisées des capacités, telles que les tests de référence à choix multiples, peuvent ne pas mesurer les capacités des systèmes d'IA dans les contextes les plus pertinents par rapport à leurs risques (par exemple lorsqu'ils sont utilisés comme aide par les humains).
- Après le développement initial, les modèles d'IA peuvent être continuellement améliorés grâce à des ajustements précis et des améliorations au niveau du temps d'inférence. Ces améliorations augmenteront les capacités contextuelles et affecteront potentiellement les risques des modèles qui sont déjà disponibles au public, et les changements seraient hors de portée des tests des développeurs du modèle de base. Il sera difficile de concevoir une politique robuste à ce type de changement continu.

1.3. Capacités dans les années à venir

INFORMATIONS CLÉS†

- Dans les mois et les années à venir, les capacités des systèmes d'IA à usage général pourraient

Les évolutions de l'IA peuvent être lentes, rapides ou extrêmement rapides. Les avis des experts et les données disponibles soutiennent chacune de ces trajectoires. Pour prendre des décisions opportunes, les décideurs politiques devront tenir compte de ces scénarios et des risques associés. Une question clé est de savoir à quelle vitesse les développeurs d'IA peuvent faire évoluer les approches existantes en utilisant encore plus de calcul et de données, et si cela suffirait à surmonter les limites des systèmes actuels, comme leur manque de fiabilité dans l'exécution de tâches longues.

- Les développeurs d'IA à usage général font progresser les domaines scientifiques, techniques et « agents »

capacités. Ces derniers mois, les modèles se sont considérablement améliorés dans les tests de raisonnement scientifique et de programmation, permettant de nouvelles applications. En outre, les développeurs d'IA déploient de gros efforts pour développer des agents d'IA polyvalents plus fiables, capables d'exécuter des tâches ou des projets plus longs sans surveillance humaine en utilisant des ordinateurs et des outils logiciels, potentiellement avec un apprentissage continu pendant le fonctionnement.

- Les outils d'IA à usage général sont de plus en plus utilisés pour accélérer

Le développement de logiciels et de matériels, y compris l'IA à usage général elle-même, est largement utilisé pour écrire plus efficacement des logiciels destinés à former et à déployer l'IA, pour aider à la conception de puces d'IA et pour générer et organiser des données de formation. L'impact de ces technologies sur le rythme des progrès a été peu étudié.

- Les améliorations récentes ont été principalement motivées par l'augmentation des capacités de calcul et de données

utilisés pour la pré-formation et en affinant les approches algorithmiques existantes. Pour les modèles de pointe, les estimations actuelles suggèrent que ces facteurs ont, ces dernières années, augmenté d'environ :

- Calcul pour pré-formation : 4x/an
- Taille de l'ensemble de données de pré-formation : 2,5 x/an
- Énergie utilisée pour alimenter les puces informatiques pendant la formation : 3x/an
- Efficacité de pré-formation algorithmique : 3x/an (incertitude plus élevée)
- Efficacité matérielle : 1,3x/an

- Il est probablement possible pour les développeurs d'IA de continuer à augmenter de manière exponentielle les ressources utilisées

Pour la formation, mais ce n'est pas garanti. Si les tendances récentes se poursuivent, d'ici la fin de 2026, les développeurs d'IA formeront des modèles en utilisant environ 100 fois plus de calcul d'entraînement que les modèles les plus gourmands en calcul de 2023, et 10 000 fois plus de calcul d'entraînement d'ici 2030. De nouvelles recherches suggèrent que ce degré de mise à l'échelle est probablement réalisable, en fonction des investissements et des décisions politiques. Cependant, il est plus probable que le rythme de mise à l'échelle actuel devienne irréalisable après les années 2020 en raison des goulots d'étranglement dans les données, la production de puces, le capital financier et l'approvisionnement énergétique local.

† Veuillez vous référer à la [mise à jour du Président](#) sur les dernières avancées en matière d'IA après la rédaction de ce rapport.

- Les chercheurs débattent de l'efficacité de l'augmentation des ressources pour la formation avec les techniques algorithmiques. Certains experts sont sceptiques quant à savoir si l'augmentation des ressources de formation suffirait à surmonter les limites des systèmes actuels, tandis que d'autres s'attendent à ce que cela continue d'être l'ingrédient clé des avancées futures.
- Les développeurs d'IA ont récemment adopté une mise à l'échelle supplémentaire potentiellement plus efficace approche. Les modèles peuvent être formés pour écrire des « chaînes de pensée » plus longues afin de décomposer les problèmes en étapes avant de générer des réponses, ce qui permet une mise à l'échelle des calculs pendant l'exécution plutôt que pendant la formation. Cette méthode s'est révélée prometteuse pour surmonter diverses limitations dans les tests de raisonnement scientifique et de programmation et peut fournir une voie supplémentaire si la mise à l'échelle de la formation traditionnelle produit des rendements décroissants.
- Depuis la publication du rapport intermédiaire (mai 2024), les systèmes d'IA à usage général ont devenir plus abordables à utiliser, plus utiles sur le plan pratique et plus largement adoptés. Les développeurs ont également considérablement amélioré les performances des modèles lors des tests de raisonnement mathématique et scientifique (voir [1.2. Capacités actuelles](#)).
- Les décideurs politiques sont confrontés à des défis pour surveiller et répondre aux progrès de l'IA. Principaux défis inclure le suivi quantitatif des avancées de l'IA et de leurs principaux moteurs, ainsi que la conception de cadres de gestion adaptative des risques qui activent les mesures d'atténuation uniquement lorsque les capacités (et les risques associés) augmentent.

Définitions clés

- Lois d'échelle : Relations systématiques observées entre la taille d'un modèle d'IA (ou la quantité du temps, des données ou des ressources informatiques utilisées dans la formation ou l'inférence) et de ses performances.
- Calcul : abréviation de « ressources informatiques », qui fait référence au matériel (par exemple les GPU), logiciels (par exemple, logiciels de gestion de données) et infrastructures (par exemple, centres de données) nécessaires pour former et exécuter les systèmes d'IA.
- Efficacité algorithmique (de formation) : un ensemble de mesures de l'efficacité avec laquelle un algorithme utilise des ressources informatiques permettant d'apprendre à partir de données, telles que la quantité de mémoire utilisée ou le temps nécessaire à la formation.
- Agent IA : une IA à usage général qui peut élaborer des plans pour atteindre des objectifs, exécuter des tâches de manière adaptative des tâches impliquant plusieurs étapes et des résultats incertains en cours de route, et interagissent avec son environnement - par exemple en créant des fichiers, en effectuant des actions sur le Web ou en déléguant des tâches à d'autres agents - avec peu ou pas de surveillance humaine.
- Inférence : Le processus dans lequel une IA génère des sorties basées sur une entrée donnée, appliquer les connaissances acquises lors de la formation.
- Chaîne de pensée : Un processus de raisonnement dans lequel une IA génère des étapes intermédiaires ou explications lors de la résolution d'un problème ou de la réponse à une question. Cette approche imite le raisonnement logique humain et la délibération interne, aidant le modèle à décomposer les tâches complexes en étapes séquentielles plus petites pour améliorer la précision et la transparence de ses résultats.
- Benchmark : un test ou une mesure standardisée, souvent quantitative, utilisée pour évaluer et comparer les performances des systèmes d'IA sur un ensemble fixe de tâches conçues pour représenter le monde réel usage.

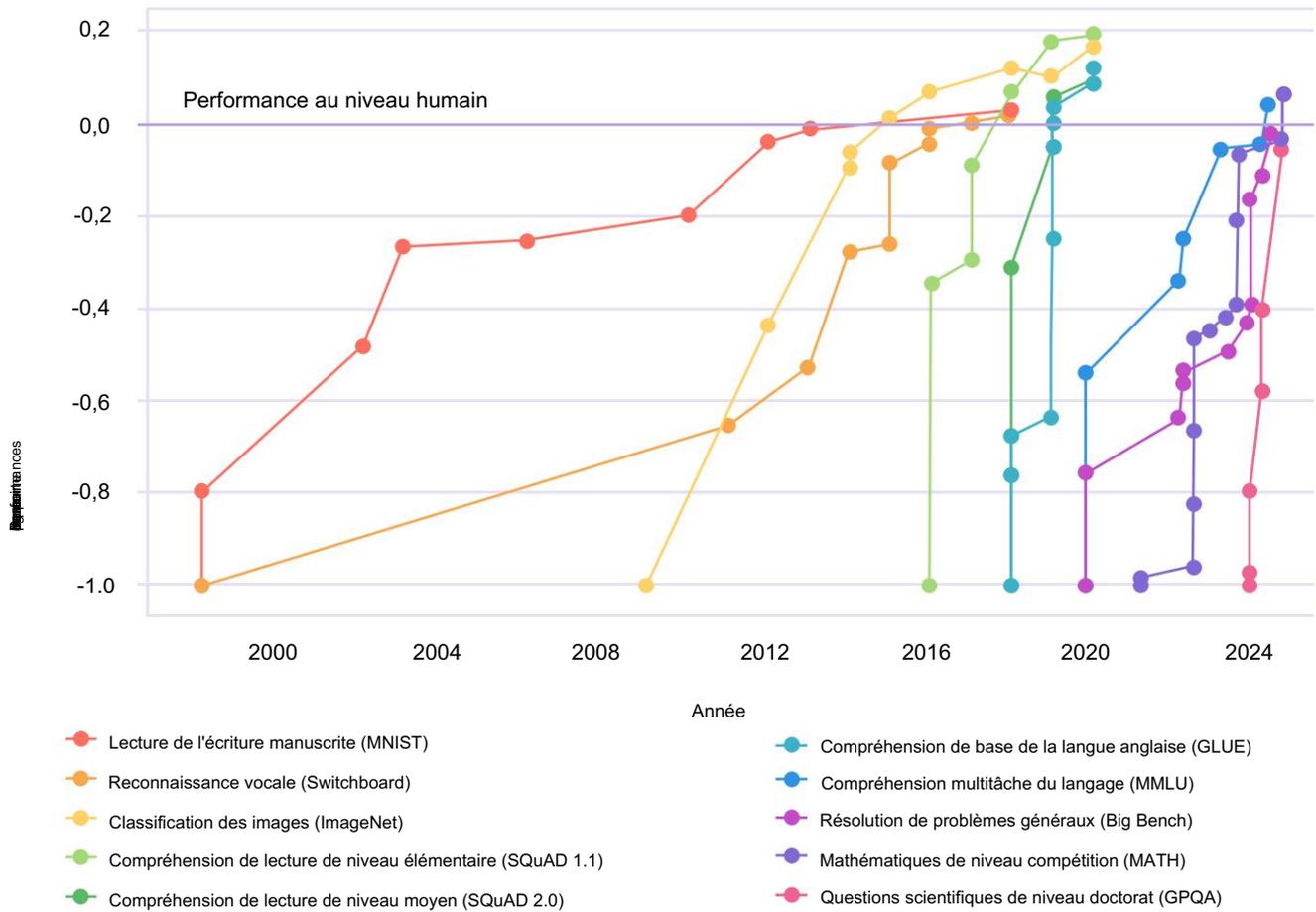
- Comportement émergent : la capacité des systèmes d'IA à agir d'une manière qui n'était pas explicitement programmés ou prévus par leurs développeurs ou utilisateurs.
- Tâches cognitives : activités qui impliquent le traitement de l'information, la résolution de problèmes, Prise de décision et pensée créative. Les exemples incluent la recherche, la rédaction et la programmation.
- Données synthétiques : données telles que du texte ou des images qui ont été générées artificiellement, par exemple par Systèmes d'IA à usage général. Les données synthétiques peuvent être utilisées pour former des systèmes d'IA, par exemple lorsque les données naturelles de haute qualité sont rares.
- Modalités : les types de données qu'un système d'IA peut recevoir avec compétence en entrée et produire en sortie, notamment du texte (langage ou code), des images, des vidéos et des actions robotiques.

1.3.1. Tendances récentes en matière de capacités d'IA à usage général

Le rythme des progrès récents en matière d'IA à usage général a été rapide, dépassant souvent les attentes des experts en IA sur la base de mesures largement utilisées. Les chercheurs évaluent les performances de l'IA à l'aide de « repères » – Français ensembles standardisés de problèmes conçus pour comparer les performances des systèmes d'IA dans un ou plusieurs domaines. Au cours de la dernière décennie, les systèmes d'IA à usage général et les systèmes d'IA antérieurs ont atteint ou dépassé les performances de niveau humain sur des tests de référence dans une grande variété de domaines, tels que le traitement du langage naturel, la vision par ordinateur, la reconnaissance vocale et les mathématiques (voir Figure 1.4). Par exemple, considérez le test de référence MATH (137), qui teste les compétences en résolution de problèmes mathématiques via une série de problèmes énoncés. Ces problèmes varient en difficulté, depuis de simples questions de niveau primaire jusqu'à des problèmes qui mettent au défi les lauréats de concours internationaux de mathématiques. Lorsque ce test de référence a été publié en 2021, les systèmes d'IA à usage général ont obtenu un score d'environ 5 %, mais trois ans plus tard, le modèle o1 a atteint 94,8 % (92*), ce qui correspond au score des testeurs humains experts (dans ce cas, un médaillé d'or de l'IMO). Cependant, il est souvent difficile de savoir comment des performances impressionnantes sur des tests de référence se traduisent en performances dans des tâches du monde réel, comme indiqué ci-dessous (138).

Les systèmes d'IA sont devenus beaucoup plus rentables à exploiter, les prix de fonctionnement des systèmes d'IA à un niveau de capacité donné diminuant de plusieurs ordres de grandeur. Par exemple, en 2022, il en coûtait environ 25 dollars aux utilisateurs pour générer un million de mots avec GPT-3, mais en 2023, ce coût est tombé à près de 1 dollar avec l'équivalent en performances Llama 2 7B (voir la figure 1.5). Ces baisses de prix découlent en partie des avancées technologiques, telles que les améliorations matérielles qui permettent d'effectuer davantage de calculs au même prix (144). Les baisses de prix peuvent également se produire en raison de la diminution des marges tarifaires facturées par les entreprises, et la baisse mesurée dépend également de la référence choisie et du niveau de performance.

Comparaison des performances de l'IA et des performances humaines sur certains benchmarks



Français : Figure 1.4 : Les performances des modèles d'IA sur divers benchmarks ont progressé rapidement entre 1998 et 2024. Notez que certains résultats antérieurs utilisaient des modèles d'IA d'apprentissage automatique qui ne sont pas des modèles à usage général. Sur certains benchmarks récents, les modèles sont passés en peu de temps d'une performance médiocre à une performance supérieure à celle de sujets humains qui sont souvent des experts. Notez que les premiers résultats de ce graphique utilisaient des modèles d'IA d'apprentissage automatique qui ne sont pas des modèles à usage général. Sources : Kiela et al., 2021 (139) (pour MNIST, Switchboard, ImageNet, SQuAD 1.1, 2 et GLUE). Les données pour MMLU, Big Bench, GPQA proviennent des articles pertinents (3*, 5*, 92*, 140, 141, 142, 143*).

Depuis la publication du rapport intermédiaire, les recherches sur l'amélioration des capacités de l'IA à usage général ont commencé à se concentrer sur de nouvelles directions, tandis que les efforts pour augmenter les ressources de formation se poursuivent. Par exemple, l'une des directions consiste à améliorer l'autonomie des systèmes d'IA à usage général, en produisant des agents d'IA qui agissent et planifient dans la poursuite d'objectifs (150) (voir [1.2. Capacités actuelles](#) et [3.2.1. Défis techniques pour la gestion des risques et l'élaboration des politiques](#)). Une autre direction consiste à utiliser plusieurs copies de modèles ensemble pour accomplir de nouvelles tâches (151*).

Les modèles linguistiques sont proposés à moindre coût, générant plus de mots par dollar

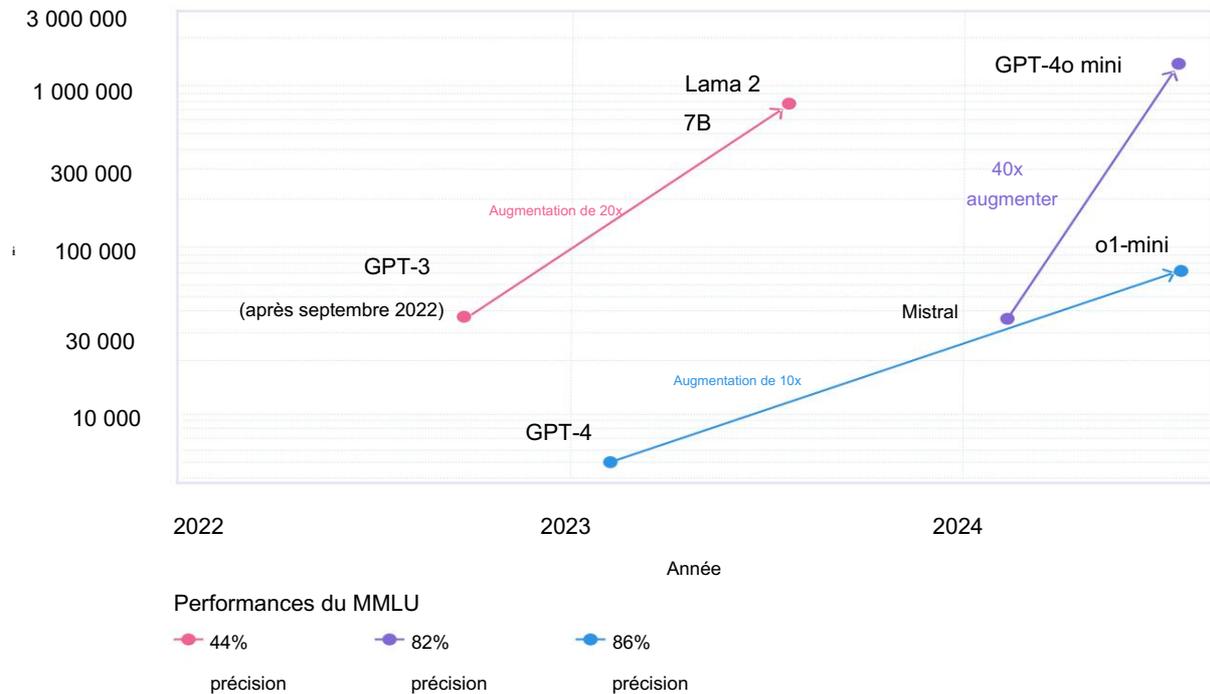


Figure 1.5 : Ce graphique montre comment les modèles linguistiques à usage général sont devenus nettement plus rentables à utiliser, mesurés par le nombre de mots générés par dollar tout en maintenant un niveau de performance donné sur le benchmark MMLU. La version de GPT-3 175B publiée après septembre 2022 et Llama 2 7B atteignent tous deux un score d'environ 44 % de précision (48*, 145*), tandis que Mistral Large et GPT-4o mini atteignent environ 82 % (12*, 146*). Le GPT-4 original de mars 2023 et le o1-mini récemment publié obtiennent tous deux un score d'environ 86 % sur MMLU (92*, 147*). Notez que ce graphique est principalement à des fins d'illustration, car les prix déclarés et les performances MMLU dépendent des méthodes d'évaluation.

Français De plus, o1-mini écrit ce que l'on appelle des « chaînes de pensée » auxquelles les utilisateurs ne peuvent pas accéder avant de produire une réponse finale, donc dans la pratique le nombre de mots accessibles générés par dollar est probablement inférieur à celui représenté dans la figure. Sources : Chung et al., 2022 (145*) et Touvron et al., 2023 (48*) (pour GPT-3 175B et Llama 2 7B) ; Mistral AI, 2024 (12*) et OpenAI, 2024f (146*) (pour Mistral Large et GPT-4o mini) ; Open AI, 2024g (92*) et OpenAI et al., 2024 (147*) (pour GPT-4 et o1-mini) ; OpenAI, 2024d (148*) et Together Pricing, 2023 (149*) (pour les données de tarification).

De nouvelles données suggèrent que la mise à l'échelle des capacités de calcul et des données de formation aux taux actuels est techniquement possible jusqu'en 2030 au moins. Au cours de la dernière décennie, la capacité de calcul de formation des modèles de pointe a augmenté d'environ 4 fois par an. Si cette tendance se poursuit, les systèmes seront formés avec environ 100 fois plus de puissance de calcul que GPT-4 d'ici la fin de 2026, et jusqu'à environ 10 000 fois d'ici la fin de la décennie (152). Cependant, on ne sait pas clairement comment cela se traduira par des capacités améliorées, et si les retours économiques seront suffisamment importants pour justifier les dépenses liées à des niveaux de mise à l'échelle aussi massifs.

1.3.2. Les limites des systèmes actuels peuvent-elles être résolues en faisant évoluer, en perfectionnant et en combinant les approches existantes ?

Les systèmes d'IA à usage général actuels disposent d'un ensemble inégal de capacités et présentent encore de nombreuses limitations

Français Les humains et les systèmes d'IA à usage général présentent des forces et des faiblesses distinctes, ce qui rend les comparaisons difficiles. Il est tentant de comparer les capacités cognitives des humains et des systèmes d'IA, par exemple parce que cela permet de savoir quelles tâches économiques pourraient être particulièrement fortement impactées par l'utilisation de l'IA. Cependant, les systèmes d'IA à usage général actuels affichent souvent des performances inégales, excellant dans certains domaines tout en étant en difficulté dans d'autres (153), ce qui rend les comparaisons trop générales moins significatives. Si l'IA à usage général surpasse désormais les humains sur certains critères, certains scientifiques affirment qu'elle ne dispose toujours pas de la compréhension conceptuelle approfondie et des capacités de raisonnement abstrait des humains (153). Les systèmes d'IA à usage général peuvent remplacer les humains dans certaines activités, tandis que dans d'autres, les forces et les faiblesses distinctes des systèmes d'IA et des humains se combinent pour produire des collaborations fructueuses (voir [2.3.1. Risques liés au marché du travail](#)).

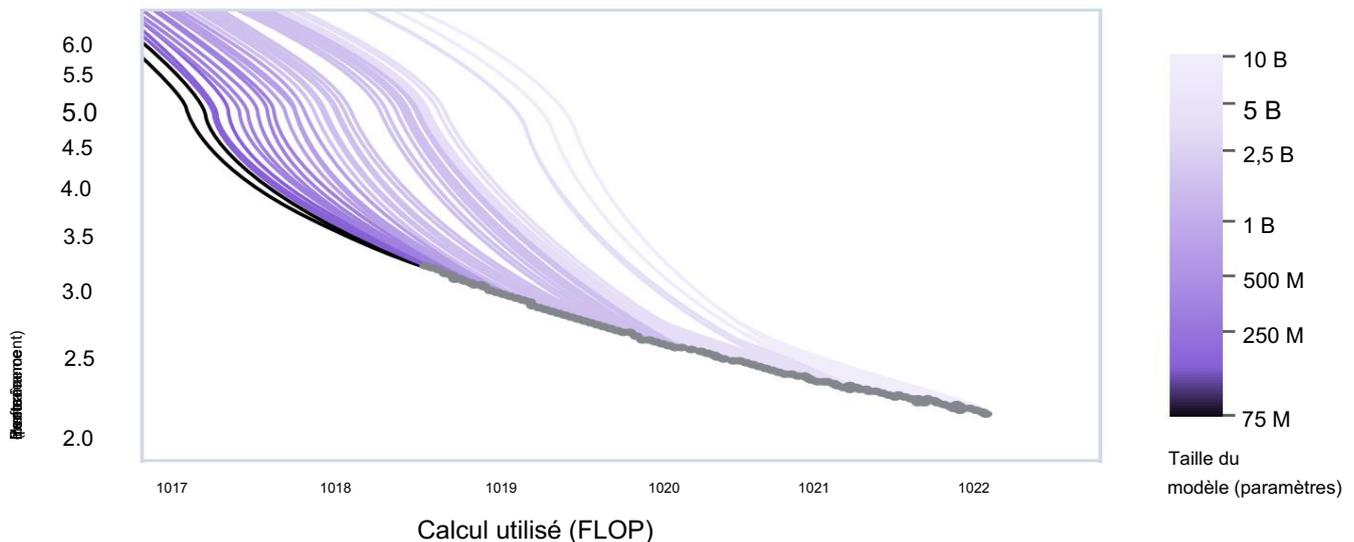
Les systèmes d'IA à usage général actuels sont sujets à certaines défaillances que les humains ne connaissent pas (154, 155). Certains travaux suggèrent que le raisonnement de l'IA à usage général peut avoir du mal à faire face à des scénarios nouveaux et est trop influencé par des similitudes superficielles (110*, 153). Il a également été démontré que les systèmes d'IA à usage général échouent parfois à raisonner sur des tâches apparemment simples. Par exemple, un modèle formé sur des données comprenant l'affirmation « Olaf Scholz était le neuvième chancelier d'Allemagne » ne sera pas toujours en mesure de répondre à la question « Qui était le neuvième chancelier d'Allemagne ? » (154). En outre, il existe des preuves que les systèmes d'IA à usage général peuvent être amenés à s'écarter de leurs protections habituelles par des entrées absurdes, alors que les humains reconnaîtraient ces invites (voir [3.4.1. Former des modèles plus fiables](#)). Les limites des systèmes actuels sont examinées plus en détail dans [1.2. Capacités actuelles](#).

Les approches de formation de l'IA existantes étendront probablement les capacités du modèle, mais le degré d'amélioration et son importance dans le monde réel font l'objet de nombreux débats.

Les données suggèrent qu'une mise à l'échelle supplémentaire des ressources augmentera les capacités globales de l'IA. Les chercheurs ont découvert des « lois d'échelle » empiriques (voir Figure 1.6), qui sont des relations mathématiques qui quantifient la relation entre les entrées du processus de formation de l'IA (telles que les quantités de données et de calcul) et les capacités du modèle sur des tâches de performance générales telles que la prédiction du mot suivant (156*, 157*). Ces études démontrent que les performances des modèles d'IA ont tendance à s'améliorer avec l'augmentation des ressources de calcul dans un éventail de domaines, notamment la vision par ordinateur (158*, 159), la modélisation du langage (156*, 157*) et les jeux

jouer (160*). Bien que de nombreuses mesures de performance ne testent pas directement les capacités du monde réel, il a été observé que les performances des modèles d'IA à usage général s'améliorent constamment par rapport aux grandes références qui testent de nombreuses capacités, telles que MMLU (140), à mesure que les modèles sont mis à l'échelle.

Les performances de prédiction du mot suivant s'améliorent de manière prévisible avec plus de calculs



Cependant, il n'est pas certain que la mise à l'échelle des ressources améliorera les capacités de l'IA au même rythme qu'au cours de la dernière décennie. Les lois d'échelle se sont révélées robustes, se maintenant sur une plage de multiplications par million ou par milliard des calculs d'entraînement. Cependant, ces lois d'échelle ont jusqu'à présent été dérivées d'observations empiriques, et non de principes inviolables (bien que des modèles théoriques aient été proposés pour les expliquer) (161*, 162*, 163, 164, 165). De plus, certaines lois d'échelle sont dérivées de données limitées, ce qui les rend moins fiables (41, 166*, 167, 168*, 169, 170*). Par conséquent, il n'existe aucune garantie mathématique que les lois d'échelle continueront à se maintenir à des échelles plus grandes, au-delà de la plage des données empiriques utilisées pour les établir. D'un autre côté, une décomposition des principales lois d'échelle n'a pas non plus été établie scientifiquement, malgré les reportages en cours.

Bien que les capacités globales de l'IA s'améliorent de manière prévisible avec l'échelle, il est difficile de prédire quand des capacités spécifiques apparaîtront. Il existe de nombreux exemples documentés de capacités qui apparaissent lorsque les modèles atteignent une certaine échelle, parfois soudainement, sans être explicitement programmées dans le modèle (170*, 171, 172, 173*, 174, 175). Par exemple, les LLM à une certaine échelle ont acquis la capacité d'additionner avec précision de grands nombres, lorsqu'ils sont invités à effectuer le calcul étape par étape. Certains chercheurs les définissent comme des capacités « émergentes » (171, 172, 173*, 174), indiquant qu'elles sont

Les capacités émergentes sont présentes dans les modèles de plus grande taille, mais pas dans les modèles de plus petite taille, et leur émergence est donc souvent difficile à prévoir à l'avance. D'un autre côté, des recherches récentes ont permis de faire quelques progrès dans la prédiction des capacités « émergentes » (176, 177). Un débat est en cours sur la question de savoir si les capacités peuvent être qualifiées d'« émergentes » : certaines définitions de l'émergence exigent que la capacité apparaisse soudainement ou de manière imprévisible à une certaine échelle (ce qui n'est pas toujours le cas), tandis que d'autres définitions exigent seulement que la capacité apparaisse à mesure que les modèles sont mis à l'échelle, sans être explicitement conçues pour avoir cette capacité.

On se demande dans quelle mesure les performances des tests de référence reflètent la compréhension ou l'utilité du monde réel. Les modèles d'IA ont fait des progrès rapides sur de nombreuses mesures de référence, mais ces tests de référence sont limités par rapport aux tâches du monde réel, et les experts se demandent si ces mesures évaluent efficacement les capacités véritablement générales (178, 179*). Les modèles d'IA à usage général de pointe présentent souvent des faiblesses inattendues ou un manque de robustesse sur certains tests de référence. Par exemple, ces systèmes sont moins performants sur des variantes rares ou plus difficiles de tâches qui ne sont pas observées dans les données d'entraînement (40, 110*). Certains chercheurs émettent l'hypothèse que cela est dû au fait que les systèmes s'appuient en partie ou entièrement sur la mémorisation de modèles plutôt que sur l'utilisation d'un raisonnement robuste ou d'une pensée abstraite (153, 180*). Dans certains cas, les modèles ont été formés sur les solutions de référence, ce qui a conduit à des performances de référence élevées bien que les modèles ne soient pas capables de bien performer sur la tâche dans des contextes réels (181, 182). Les modèles ont également du mal à s'adapter à des cultures qui sont moins représentées dans les données d'entraînement (183). Des problèmes comme ceux-ci soulignent la difficulté d'évaluer ce que les résultats de référence impliquent sur la capacité des modèles à appliquer de manière fiable les connaissances à des scénarios pratiques du monde réel.

Cependant, les systèmes d'IA à usage général ont parfois de bons résultats sur des tâches difficiles conçues pour tester le raisonnement, sans avoir eu la possibilité de mémoriser les solutions. En général, la présence de mémorisation constatée dans certaines études n'implique pas l'absence de processus plus avancés comme le raisonnement – il est possible que les deux existent dans des modèles différents ou au sein du même modèle. Il existe des preuves (184*, 185) que certains modèles d'IA ont généralisé leur apprentissage à des situations sur lesquelles ils n'ont pas été formés, ce qui suggère qu'ils ne se contentent pas de mémoriser des données. Certains modèles de langage à usage général (et les systèmes construits avec eux) ont obtenu de bons résultats sur des problèmes de raisonnement et de mathématiques dont les solutions ne faisaient pas partie de leurs données d'entraînement (186*). Cela s'étend à l'obtention de performances de niveau médaille aux récentes Olympiades internationales de mathématiques (187*, 188) et d'informatique (92*) et au difficile Corpus d'abstraction et de raisonnement (ARC, (189)).

Il existe un désaccord important sur la question de savoir si les développeurs d'IA peuvent atteindre un niveau d'IA largement humain pour la plupart des tâches cognitives en faisant évoluer les ressources de formation ainsi qu'en affinant et en combinant les techniques existantes. Certains soutiennent qu'une mise à l'échelle continue (éventuellement combinée à un perfectionnement et à une combinaison des approches existantes) pourrait conduire au développement de systèmes d'IA polyvalents qui fonctionnent à un niveau largement humain ou au-delà pour la plupart des tâches cognitives (190). Ce point de vue est étayé par l'observation de lois d'échelle cohérentes et par la façon dont l'augmentation de l'échelle a permis de surmonter de nombreuses limitations des premiers modèles de langage tels que GPT-1, qui pouvaient rarement générer un paragraphe de texte cohérent.

D'autres soutiennent que l'apprentissage profond présente des limites fondamentales qui ne peuvent être résolues par la seule mise à l'échelle. Ces critiques soutiennent que les systèmes actuels reposent sur la mémorisation (au moins partiellement, voir

(ci-dessus), et manquent de véritable raisonnement de bon sens (153, 191, 192), de raisonnement causal (193) ou d'une compréhension du monde physique (153, 191, 193), ainsi que d'autres limitations discutées dans 1.2. [Capacités actuelles](#). Pour remédier aux limitations actuelles, affirment-ils, il faudra peut-être des avancées conceptuelles et des innovations importantes au-delà du paradigme actuel de l'apprentissage profond et de la mise à l'échelle. Cependant, avec la découverte de o1 (2*), les chercheurs ont récemment identifié une méthode de mise à l'échelle potentiellement plus efficace qui pourrait surmonter les limitations précédentes ou servir d'alternative si les retours sur investissement de la mise à l'échelle traditionnelle diminuent de manière significative (voir 1.2. [Capacités actuelles](#)).

1.3.3. Dans quelle mesure les approches existantes seront-elles améliorées et perfectionnées dans les années à venir ?

Les ressources informatiques dédiées à la formation de l'IA ont été rapidement augmentées, et une nouvelle augmentation rapide jusqu'en 2030 semble envisageable

Les développeurs d'IA ont augmenté rapidement la capacité de formation des modèles phares, avec une croissance d'environ 4x/an. L'utilisation des capacités de calcul pour l'entraînement a connu une croissance exponentielle depuis le début des années 2010 (voir la figure 1.7), la quantité moyenne utilisée pour entraîner les modèles d'apprentissage automatique doublant environ tous les six mois (26). À titre d'illustration, les modèles d'apprentissage automatique notables (194, 195, 196) en 2010 ont été entraînés avec environ dix milliards de fois moins de puissance de calcul que les plus grands modèles en 2023 (197, 198*).

Les entreprises d'IA ont également investi davantage de ressources informatiques dans le déploiement. Cela s'explique à la fois par le fait que davantage de systèmes d'IA à usage général ont été déployés pour servir les utilisateurs (199) et par le fait que les systèmes déployés ont accès à davantage de ressources informatiques pour accroître leurs capacités. Les modèles peuvent être exécutés plus longtemps ou les résultats de plusieurs modèles peuvent être agrégés, ce qui entraîne des gains de performances qui complètent les gains liés à l'utilisation de davantage de calculs d'entraînement (80*, 92*, 93, 94*, 200*, 201, 202*, 203*, 204). Par exemple, certaines estimations indiquent qu'OpenAI a engagé 700 000 \$ par jour de coûts de déploiement en 2023 (205) et que l'exécution de l'IA représentait 60 % des émissions de CO2 de Google provenant de l'infrastructure d'apprentissage automatique en 2022 (206).

La quantité de calcul d'entraînement disponible n'a cessé de croître, principalement en raison des dépenses d'investissement importantes qui ont augmenté la quantité de puces d'IA. Depuis 2010, le matériel informatique est devenu moins cher en raison des améliorations matérielles, ce qui signifie que la quantité de puissance de calcul (calcul) que les entreprises d'IA peuvent acheter avec un dollar augmente à un rythme de 1,35 fois par an (144, 207). Cependant, la puissance de calcul totale utilisée pour entraîner des systèmes d'IA notables a augmenté d'environ 4 fois par an depuis 2010 (26), dépassant le taux d'amélioration de l'efficacité du matériel. Cela suggère que le principal moteur de la croissance du calcul d'entraînement a été les investissements visant à étendre le parc de puces d'IA, et non l'amélioration des performances des puces.

Les calculs de l'IA ont des besoins énergétiques considérables, mais les taux de croissance actuels de la consommation d'énergie de l'IA pourraient persister pendant plusieurs années. Le calcul de l'IA à l'échelle mondiale devrait nécessiter de l'électricité

consommation similaire à celle de l'Autriche ou de la Finlande d'ici 2026 (208) (voir [2.3.4 Risques pour l'environnement](#) (pour plus d'informations). Sur la base des taux de croissance actuels de la consommation d'énergie pour la formation de l'IA, les plus grandes séries de formation de l'IA en 2030 nécessiteront 1 à 5 gigawatts (GW) d'énergie. En effet, un fournisseur de calcul a récemment acheté un centre de données avec une alimentation électrique de 960 mégawatts (209). Ainsi, en fonction des décisions d'investissement et de politique, les goulots d'étranglement énergétiques n'empêcheront probablement pas le calcul d'évoluer aux taux actuels jusqu'à la fin de la décennie.

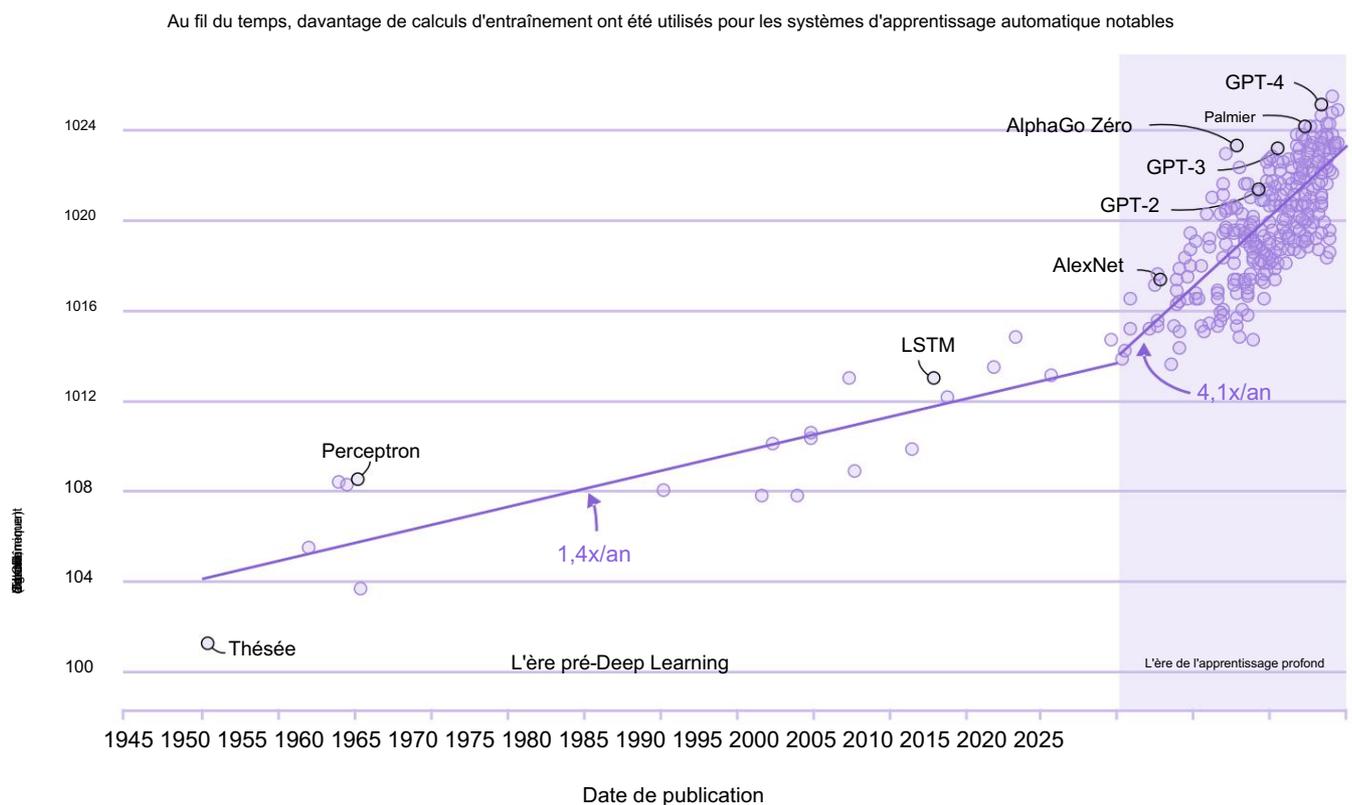


Figure 1.7 : Les développeurs d'IA ont constamment utilisé davantage de calcul pour former des modèles d'apprentissage automatique notables au fil du temps, à un rythme croissant depuis 2010 (26, 197). Le calcul est mesuré en FLOP total (opérations à virgule flottante) estimé à partir de la littérature sur l'IA — cela fait référence au nombre d'opérations de calcul effectuées pendant la formation. Les estimations devraient être précises dans un facteur de deux, ou d'un facteur de cinq pour les modèles récents non divulgués tels que GPT-4. Sources : Epoch AI, 2024 (26, 197) ; Sevilla et al., 2022 (26, 197).

La production et l'amélioration des puces d'IA sont confrontées à des défis, mais il est vraisemblable qu'ils pourront être surmontés. Il faut généralement 3 à 5 ans pour construire une usine de fabrication de puces informatiques (210, 211), et les pénuries dans la chaîne d'approvisionnement retardent parfois la production de composants de puces importants (212, 213, 214). Cependant, les grandes entreprises d'IA peuvent encore soutenir la croissance du calcul à court terme en capturant une grande partie du stock de puces d'IA. Par exemple, une étude estime que la part des puces d'IA des centres de données du monde entier détenues par une seule entreprise d'IA à tout moment se situe entre 10 % et 40 % (215). De plus, une analyse des tendances et des possibilités techniques existantes dans la production de puces suggère qu'il est possible de former des systèmes d'IA avec 100 000 fois plus de calcul d'entraînement que GPT-4 (le principal modèle de langage de 2023) d'ici 2030. Cela est suffisant pour soutenir les taux de croissance existants dans le calcul d'entraînement, qui impliquent une augmentation totale de 10 000 fois sur la même période (215).

Les contraintes de production sont importantes, mais il est peu probable qu'elles empêchent une mise à l'échelle supplémentaire des plus grands modèles aux rythmes actuels jusqu'en 2030 si les investissements sont maintenus (voir figure 1.8).

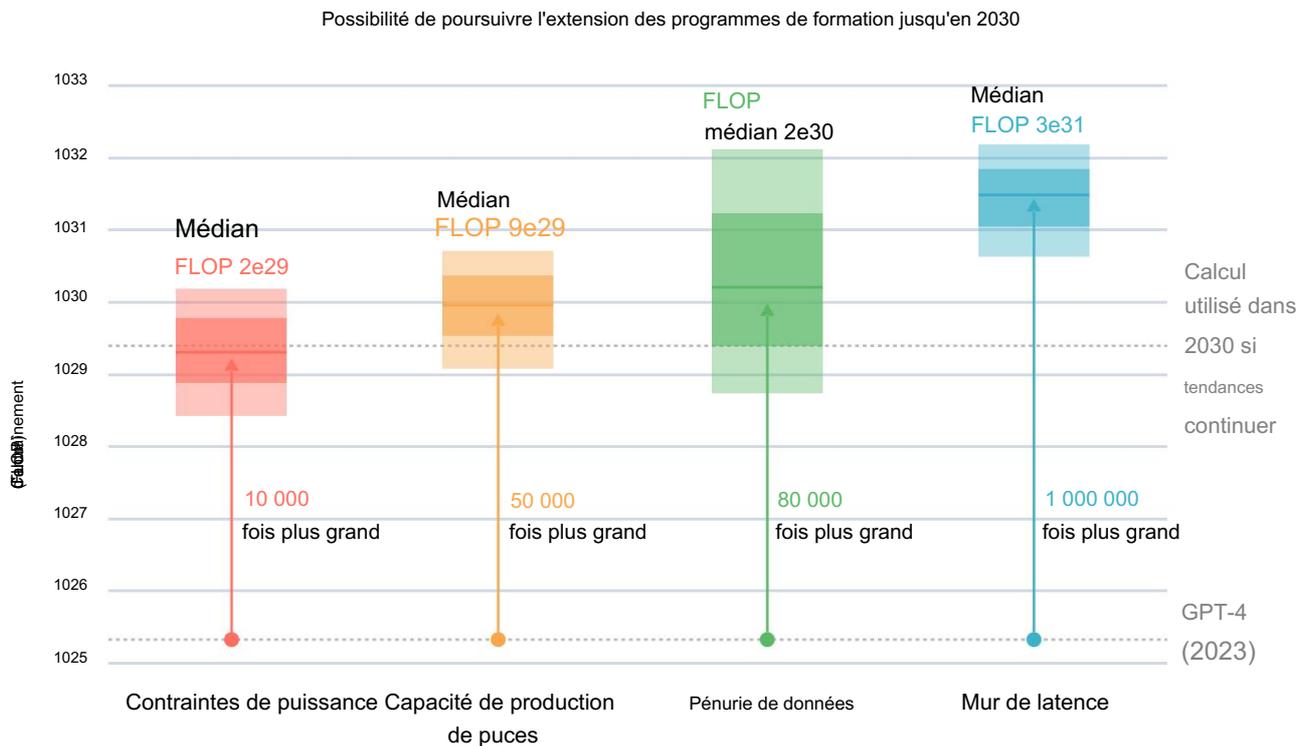


Figure 1.8 : Quatre contraintes physiques à la formation de modèles d'IA à usage général utilisant davantage de calcul d'ici 2030. Il existe de nombreux ordres de grandeur d'incertitude dans les estimations globales, mais les exécutions de formation utilisant 10 000 fois plus de calcul que GPT-4 (sorti en 2023), ce qui est conforme à la tendance actuelle, semblent techniquement réalisables sur la base de ces estimations. Source : Sevilla et al. 2024 (215).

La formation des systèmes d'IA sur un très grand nombre de puces d'IA est difficile, ce qui peut empêcher des cycles d'entraînement extrêmement importants. Par exemple, certaines estimations suggèrent que des cycles d'entraînement 10 000 à 10 millions de fois plus importants que ceux de GPT-4 seront impossibles en raison des contraintes sur la quantité d'informations pouvant être transférées entre les puces et des limites sur le temps de traitement des données (215, 216). Si ces estimations sont correctes, ces goulots d'étranglement limiteront la capacité des développeurs d'IA à augmenter les capacités de calcul d'entraînement aux taux actuels au cours de la prochaine décennie. Cependant, il est possible que de nouvelles techniques ou des solutions de contournement simples permettent des cycles d'entraînement beaucoup plus importants.

Il existe probablement suffisamment de données de pré-formation pour une mise à l'échelle jusqu'en 2030, mais les projections sont très incertaines après cette date.

Les pénuries de données constituent un goulot d'étranglement plausible pour la mise à l'échelle continue du pré-entraînement du modèle linguistique. Depuis 2010, les besoins en données pour la préformation des systèmes d'IA à usage général ont augmenté d'environ 10 fois tous les trois ans (197). Par exemple, un modèle de pointe en 2017 a été formé avec quelques milliards de mots, tandis que les modèles à usage général de pointe en 2023 ont été formés avec plusieurs milliers de milliards (217*, 218*). Une grande partie de cette croissance a été possible grâce à la disponibilité des données sur Internet, mais les taux de croissance

La demande de données semble suffisamment rapide pour épuiser les données textuelles générées par l'homme sur Internet d'ici 2030 (219, 220). Ces défis sont exacerbés par les problèmes de droits d'auteur sur les données, car il pourrait devenir illégal pour les entreprises d'IA de former l'IA sur certains types de données (voir [2.3.6. Risques de violation des droits d'auteur](#)).

Le degré de rareté des données est spécifique au domaine et à l'acteur. Dans certains domaines, la collecte de données peut être considérablement augmentée, comme dans la robotique à usage général, où les systèmes collectent des données pendant le déploiement (221*).

L'approvisionnement en données de différentes modalités pourrait contribuer à soutenir la mise à l'échelle des données. Les systèmes d'IA à usage général sont de plus en plus formés sur des données multimodales, combinant des informations textuelles, visuelles, auditives ou biologiques (59*, 222, 223, 224*). Plusieurs études suggèrent que cela augmentera la quantité de données de formation disponibles pour les modèles et dotera les modèles de nouvelles capacités, telles que la capacité d'analyser des documents contenant à la fois du texte et des graphiques (4*, 50*, 147*). Les estimations les plus complètes suggèrent qu'il existe suffisamment de données multimodales pour prendre en charge des exécutions de formation mille à dix millions de fois plus importantes que celles de GPT-4 en termes de taille de calcul, ce qui nécessite environ dix fois plus de données (215, 225). Cependant, ces estimations sont très incertaines, car il est difficile d'évaluer dans quelle mesure la formation sur une modalité de données affecte les performances sur une autre modalité.

Les données synthétiques générées par les machines pourraient réduire considérablement les goulets d'étranglement des données, mais les preuves de leur utilité sont mitigées. Les ensembles de données d'entraînement peuvent également être complétés par des sorties d'IA « synthétiques » à usage général, qui peuvent être utiles lorsque les données réelles sont limitées (226*, 227) ou pour améliorer la généralisation du modèle (227, 228). Cependant, certains soutiennent que l'entraînement naïf sur des sorties d'IA à usage général dégrade les performances ou a des rendements rapidement décroissants (229, 230, 231, 232, 233, 234*, 235, 236). D'autres soutiennent que ces problèmes peuvent être contournés grâce à de meilleures techniques de formation, comme l'intégration de données « naturelles » (229, 231, 235, 237*, 238), l'amélioration de la qualité des données (par exemple) en utilisant un modèle pour évaluer sa qualité (226*, 239*, 240, 241) et la formation sur des exemples négatifs (c'est-à-dire en enseignant à l'IA ce qu'elle ne doit pas faire) (242*). Les modèles phares récents tels que Llama 3 ont fait un usage substantiel de données synthétiques au cours de plusieurs étapes de formation (37*). Les récentes améliorations du modèle o1 dans les tests de raisonnement et de programmation ont été obtenues en grande partie en apprenant à partir de ses propres « chaînes de pensée » auto-générées - en analysant les chemins de raisonnement qui ont conduit au succès ou à l'échec (2*).

La plupart des succès obtenus avec les données synthétiques se limitent à certains domaines. L'apprentissage des données synthétiques peut être très efficace dans les domaines où les résultats des modèles peuvent être formellement vérifiés, comme les mathématiques et la programmation (187*, 188, 243, 244*). Cependant, il n'est pas encore certain que les méthodes d'apprentissage des données synthétiques soient efficaces dans les domaines où les résultats ne peuvent pas être facilement vérifiés. La recherche médicale en est un exemple : les données doivent souvent être vérifiées en effectuant des expériences qui durent des mois, voire des années.

1.3.4. Dans quelle mesure les capacités de l'IA seront-elles améliorées grâce à l'invention ou au perfectionnement d'algorithmes ?

Les techniques et méthodes de formation existantes pour l'IA à usage général ont été constamment améliorées et affinées

Les améliorations algorithmiques permettent de former des modèles d'IA à usage général avec moins de ressources. Les techniques et méthodes de formation qui sous-tendent les modèles d'IA polyvalents les plus performants se sont constamment et régulièrement améliorées au fil du temps. L'efficacité de calcul des techniques d'IA pour la formation a été multipliée par 10 environ tous les 2 à 5 ans dans des domaines clés tels que la classification d'images, les jeux et la modélisation du langage (245*, 246). Par exemple, la quantité de calcul nécessaire pour entraîner un modèle à atteindre un niveau de performance défini en matière de classification d'images a diminué de 44 fois entre 2012 et 2019, ce qui signifie que l'efficacité a doublé tous les 16 mois.

Les systèmes d'IA de jeu nécessitent deux fois moins d'exemples d'entraînement tous les 5 à 20 mois (247). Dans la modélisation du langage, le calcul nécessaire pour atteindre un niveau de performance fixe a diminué de moitié environ tous les huit mois en moyenne depuis 2012 (246). Cela correspond à une amélioration de l'efficacité de l'entraînement algorithmique de 3 fois par an, soit une amélioration totale d'environ 27 fois d'ici fin 2026. Ces avancées ont permis aux chercheurs et aux entreprises d'IA à usage général de développer des modèles plus performants au fil du temps dans le cadre d'un budget matériel limité.

Les innovations algorithmiques se produisent également dans d'autres dimensions, mais elles sont moins bien mesurées. Par exemple, de nouvelles techniques ont permis aux systèmes d'IA à usage général de traiter de plus grandes quantités d'informations contextuelles pour chaque requête adressée au système d'IA (248*). Certaines innovations algorithmiques contribuent également à accroître les performances, permettent aux systèmes d'IA à usage général d'utiliser des outils (22*) et de mieux exploiter le calcul lors du déploiement (94*). Ces capacités varient selon les différentes dimensions, leurs taux d'amélioration sont difficiles à mesurer et elles sont souvent moins bien comprises.

Les améliorations apportées après la préformation peuvent être utilisées pour améliorer considérablement les capacités des modèles d'IA à usage général à faible coût. Il existe un nombre croissant de travaux sur les innovations algorithmiques après la formation initiale, telles que l'amélioration du réglage fin, l'accès des modèles aux outils logiciels et la structuration des résultats des modèles pour les tâches de raisonnement (voir 1.2. [Capacités actuelles](#)). Cela signifie qu'un large éventail d'acteurs, y compris les acteurs à faibles ressources, pourraient utiliser des améliorations (parfois appelées « améliorations post-formation ») pour faire progresser les capacités d'IA à usage général - un facteur important dont la gouvernance doit tenir compte.

Progrès des capacités issues de l'application des systèmes d'IA au développement de l'IA

Les systèmes d'IA à usage général sont de plus en plus déployés pour automatiser et accélérer la recherche et le développement de l'IA, et leurs effets sur le rythme des progrès sont sous-étudiés. Des systèmes d'IA à usage restreint ont déjà été utilisés pour développer et améliorer des algorithmes (249, 250) et concevoir les dernières puces d'IA

(251). Les LLM récents sont largement utilisés dans les domaines liés à la R&D en IA, notamment dans la programmation (55), la génération et l'optimisation des invites et des paramètres de formation (252, 253, 254, 255), la supervision en remplaçant les données de rétroaction humaine (256*) et la sélection de données de formation de haute qualité (257*). Des prototypes récents ont également utilisé les LLM pour proposer de nouvelles idées de recherche (258*). Un système basé sur le LLM récemment publié a obtenu des résultats compétitifs avec des équipes humaines typiques dans des compétitions d'ingénierie en IA du monde réel (125*). Une étude récente comparant les systèmes d'IA à des ingénieurs humains experts a révélé que des agents d'IA soigneusement réglés, construits sur des modèles de pointe, obtenaient des résultats comparables à ceux des humains sur des tâches d'ingénierie de recherche en IA qui prennent généralement huit heures aux ingénieurs (259). Les agents IA ont montré de meilleures performances que les humains sur des tâches de moins de huit heures, mais ont pris du retard sur des tâches plus longues, suivant un schéma typique observé dans les performances de l'IA. Les tâches d'ingénierie de l'IA consomment la plus grande partie du temps de recherche et développement de l'IA, ce qui rend l'application de l'IA à ces tâches particulièrement importante (260). À mesure que les capacités des systèmes d'IA à usage général progressent, leur effet global sur le progrès algorithmique et l'ingénierie de l'IA nécessitera davantage de recherches pour être compris.

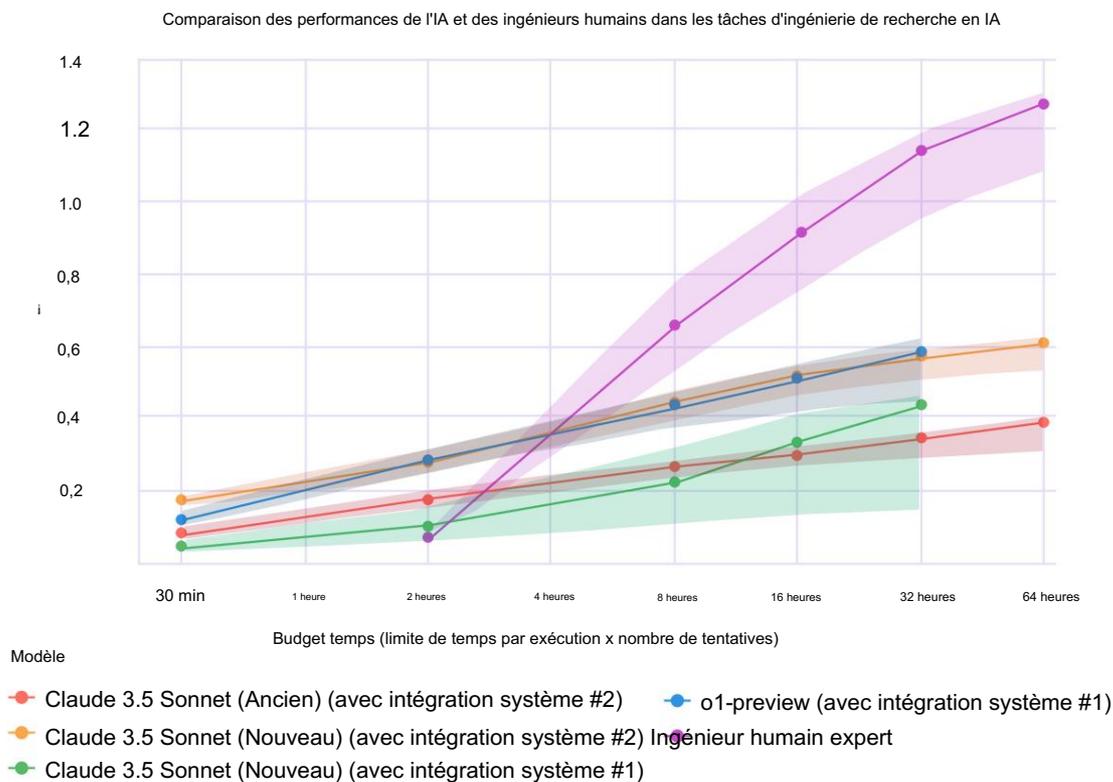


Figure 1.9 : Dans une série d'expériences, les agents d'IA basés sur le LLM publiés en 2024 ont obtenu de meilleurs résultats que les ingénieurs humains experts dans les tâches d'ingénierie de recherche en IA à durée indéterminée lorsqu'ils disposaient de deux heures ou moins pour effectuer le travail. À l'inverse, les experts humains ont obtenu de meilleurs résultats lorsqu'ils disposaient de huit heures ou plus. Différentes « intégrations de systèmes » font référence à différentes manières d'utiliser le même modèle, ce qui peut entraîner des performances variables. Les régions ombrées correspondent à des intervalles de confiance à 95 %. Source : Wijk et al., 2024 (259).

L'invention de nouvelles approches conduira-t-elle à des progrès rapides dans les années à venir ?

Les améliorations soudaines, importantes et généralisées des algorithmes d'IA sont rares mais ne peuvent être exclues.

Les percées conceptuelles fondamentales sont rares et difficiles à prévoir, car les données à ce sujet sont relativement rares. De tels événements rares ne peuvent pas être facilement prévus en extrapolant les tendances passées. Au mieux, les modèles statistiques qui analysent les améliorations passées des tests de référence de l'IA trouvent des preuves suggérant que des améliorations soudaines et importantes des performances sont peu probables, mais ne peuvent être exclues (261*).

Le corpus de preuves dans ce domaine est très limité et une incertitude substantielle demeure.

Même si les développeurs parviennent à des avancées conceptuelles fondamentales dans les algorithmes, elles ne se traduiront pas immédiatement par des améliorations de capacités importantes. Par exemple, une étude a révélé que certaines innovations algorithmiques ont des effets plus prononcés à des échelles plus grandes qu'à des échelles plus petites de calcul d'entraînement (262*), ce qui rend difficile l'observation d'améliorations dans les petites expériences. Les innovations algorithmiques doivent également être optimisées pour fonctionner correctement avec le matériel existant, ou pour être intégrées à l'infrastructure existante ou aux conventions des développeurs (263, 264*, 265*). Ces éléments constituent des obstacles à la mise en œuvre à grande échelle, donc si une avancée conceptuelle majeure est nécessaire pour surmonter les limites de l'IA à usage général actuelle, cela pourrait prendre de nombreuses années.

Défis politiques

À mesure que ces tendances techniques se poursuivent, les décideurs politiques sont confrontés à de nouveaux défis pour répondre aux impacts sociétaux de l'IA à usage général.

L'un des défis auxquels sont confrontés les décideurs politiques est la disponibilité limitée de données d'évaluation de qualité sur les capacités de l'IA à usage général. Par exemple, l'une des principales lacunes des référentiels actuels est qu'ils ne représentent pas toujours avec précision les capacités du monde réel.

Par conséquent, les efforts visant à élaborer des critères de référence plus exigeants et à mettre en place des équipes spécialisées dans l'évaluation des capacités des modèles se sont accrus (266*, 267, 268, 269).

Les problèmes de qualité des données sont encore aggravés par la quantité limitée de données, ce qui signifie que certaines estimations du taux de progrès de l'IA (par exemple pour l'efficacité algorithmique) sont très incertaines.

L'un des principaux défis consiste à gérer l'incertitude entourant la trajectoire des capacités futures. Différentes capacités d'IA à usage général pourraient avoir des répercussions très différentes sur les impacts sociétaux et les politiques en matière d'IA. Par exemple, les meilleures estimations du taux de progrès algorithmique sont très incertaines, mais le taux spécifique a des implications importantes pour les approches politiques qui mettent l'accent sur le suivi de l'utilisation des ressources informatiques pour la formation (270). Dans l'ensemble, il existe une grande incertitude quant aux capacités futures de l'IA, et des travaux supplémentaires sur le suivi des progrès de l'IA (par exemple avec des repères améliorés), ainsi que sur l'anticipation des progrès futurs, seraient utiles.

2. Risques

2.1. Risques liés à une utilisation malveillante

2.1.1. Dommages causés aux personnes par des contenus falsifiés

INFORMATIONS CLÉS

- Les acteurs malveillants peuvent utiliser l'IA à usage général pour générer du faux contenu nuisible
Les contenus frauduleux peuvent être utilisés de manière ciblée pour escroquer des personnes, faire de l'extorsion, manipuler psychologiquement, générer des images intimes non consensuelles (NCII) et du matériel pédopornographique (CSAM), ou saboter de manière ciblée des individus et des organisations.
- Cependant, les preuves scientifiques sur ces utilisations sont limitées. Des rapports anecdotiques sur les dommages causés par les faux contenus générés par l'IA sont courants, mais il n'existe pas de statistiques fiables sur la fréquence et l'impact de ces incidents. Il est donc difficile de faire des déclarations précises sur les dommages causés par les faux contenus générés par l'IA à usage général.
- Au cours des derniers mois, des progrès limités ont été réalisés dans la détermination scientifique de l'étendue de
Le problème. Depuis la publication du rapport intermédiaire (mai 2024), de nouvelles preuves suggèrent une augmentation significative de la prévalence du contenu deepfake généré par l'IA en ligne. Dans l'ensemble, les données fiables sur l'ampleur réelle du problème restent limitées.
- Plusieurs techniques d'atténuation existent, mais elles présentent toutes de sérieuses limites. Détection
Les techniques d'authentification des médias peuvent parfois aider à identifier le contenu généré par l'IA à usage général, mais des défis fondamentaux demeurent. Les techniques d'authentification des médias telles que les filigranes peuvent fournir une ligne de défense supplémentaire, mais des acteurs moyennement qualifiés peuvent généralement les supprimer.

Définitions clés

- Contenu factice généré par l'IA : contenu audio, texte ou visuel, produit par l'IA générative, qui dépeint des personnes ou des événements d'une manière qui diffère de la réalité d'une manière malveillante ou trompeuse, par exemple en montrant des personnes faisant des choses qu'elles n'ont pas faites, en disant des choses qu'elles n'ont pas dites, en changeant le lieu d'événements réels ou décrivant des événements qui n'ont pas eu lieu.
- Deepfake : un type de faux contenu généré par l'IA, composé de contenu audio ou visuel, qui déforme les faits et présente des personnes réelles comme faisant ou disant quelque chose qu'elles n'ont pas réellement fait ou dit.

Les acteurs malveillants peuvent utiliser à mauvais escient les faux contenus générés par l'IA pour extorquer, escroquer, manipuler psychologiquement ou saboter des individus ou des organisations ciblées (voir le tableau 2.1) (271). Cela menace les droits humains universels, par exemple le droit de ne pas porter atteinte à son honneur et à sa réputation (272). Cette section se concentre sur les préjudices causés aux individus par les faux contenus générés par l'IA. Les impacts potentiels des campagnes d'influence générées et médiatisées par l'IA au niveau sociétal sont abordés dans la [section 2.1.2. Manipulation de l'opinion publique.](#)

Escroqueries / fraudes	Utiliser l'IA pour générer du contenu tel qu'un clip audio imitant la voix d'une victime afin, par exemple, d'autoriser une transaction financière.
Chantage / extorsion	Générer du faux contenu d'une personne, comme des images intimes, sans son consentement et menacer de les divulguer à moins que les exigences financières ne soient satisfaites.
Sabotage	Générer du faux contenu qui présente un individu se livrant à des activités compromettantes, telles qu'une activité sexuelle ou la consommation de drogues, puis diffuser ce contenu afin de ternir la réputation d'une personne, de nuire à sa carrière et/ou de la forcer à se désengager des activités publiques (par exemple en politique, en journalisme ou dans le divertissement).
Violence	Générer des représentations néfastes d'un individu dans le but premier de l'abuser et de lui causer un traumatisme psychologique. Les victimes sont souvent des enfants.

Tableau 2.1 : Le faux contenu généré par l'IA a été utilisé pour causer différents types de préjudices à des individus, notamment par le biais d'escroqueries, de chantage, de sabotage et de violences psychologiques.

Français L'une des principales lacunes en matière de preuves concernant les dommages causés aux individus par les faux contenus est le manque de statistiques complètes et fiables sur les dommages susmentionnés, ce qui rend difficile l'évaluation précise de leur fréquence et de leur gravité. De nombreux experts pensent que les faux contenus générés artificiellement, et en particulier les contenus sexuels, sont en augmentation, mais la plupart des récits de tels cas restent anecdotiques. Les principales lacunes en matière de preuves empiriques concernent la prévalence des fraudes financières deepfake et les cas d'extorsion et de sabotage. La réticence à signaler les cas peut contribuer à ces difficultés à comprendre l'impact complet des contenus générés par l'IA destinés à nuire aux individus. Par exemple, les institutions hésitent souvent à révéler leurs problèmes de fraude alimentée par l'IA. De même, les personnes attaquées avec des documents compromettants générés par l'IA sur elles-mêmes peuvent garder le silence par gêne et pour éviter de nouveaux dommages (273).

Les criminels peuvent exploiter les faux contenus générés par l'IA pour se faire passer pour des figures d'autorité ou des personnes de confiance afin de commettre une fraude financière. Il existe de nombreux cas dans lesquels les criminels ont utilisé des clips audio et vidéo générés artificiellement pour inciter les individus à transférer de l'argent. Par exemple, les attaques de phishing peuvent exploiter un faux contenu généré par l'IA pour rendre les messages, les appels ou les vidéos frauduleux plus convaincants et efficaces, dans le but d'obtenir des informations sensibles ou de l'argent en se faisant passer pour une entité de confiance (273, 274). Les incidents vont des cas de fraude très médiatisés où des cadres de banque ont été persuadés de transférer des millions de dollars, à des individus ordinaires trompés pour transférer de plus petites sommes à des proches (soi-disant) dans le besoin. Le faux contenu généré par l'IA peut également être utilisé pour le vol d'identité, par lequel la voix ou l'image d'une victime est utilisée pour autoriser des virements bancaires ou pour ouvrir de nouveaux comptes bancaires au nom d'une victime. Alternativement, un faux contenu peut également être utilisé pour tromper les administrateurs système afin qu'ils partagent des informations de mot de passe et de nom d'utilisateur qui peuvent faciliter le vol d'identité à une date ultérieure (275).

Les faux contenus générés par l'IA peuvent également être utilisés à des fins de chantage à des fins d'extorsion. Dans de tels cas, les criminels exigent de l'argent, des secrets d'affaires ou des images ou vidéos de nus, en utilisant comme levier un contenu compromettant et réaliste généré par l'IA (276). Différents types de faux contenus générés par l'IA – allant de la vidéo, aux clones vocaux, aux images, etc. – peuvent varier en termes de réalisme et d'efficacité (277).

Les faux contenus peuvent contenir toute activité compromettante ou portant atteinte à la réputation, mais ils ont fait l'objet d'une attention particulière dans les cas de pornographie deepfake, où l'IA à usage général est utilisée pour créer des représentations audiovisuelles pornographiques ou autres représentations audiovisuelles intimes d'individus sans leur consentement (278, 279, 280). Ce contenu est ensuite utilisé pour extorquer une rançon aux victimes – en exigeant de l'argent pour empêcher la diffusion des images – ou pour obtenir le respect d'autres exigences, comme la fourniture de contenu illicite supplémentaire.

De tels contenus compromettants peuvent également être utilisés pour saboter la vie personnelle et professionnelle d'individus, violant ainsi le droit de l'homme contre les atteintes à l'honneur et à la réputation (272).

Images et vidéos fausses et compromettantes – telles que des images d'athlètes professionnels prenant des drogues –

Les cas de fraudes informatiques ont parfois entraîné des atteintes à la réputation, entraînant des pertes d'opportunités et des ruptures d'accords commerciaux (271). La possibilité de devenir l'objet de contenus deepfake préjudiciables et la menace associée de dommages à la réputation et de violences psychologiques envers soi-même et sa famille peuvent amener les gens à se désengager d'activités visibles publiquement telles que la politique et le journalisme, même lorsqu'ils n'ont pas été directement ciblés (281). Cependant, la gravité de cet « effet de silence » est difficile à estimer avec précision, car les preuves à ce stade sont en grande partie anecdotiques.

Français Les abus utilisant de faux contenus pornographiques ou intimes ciblent en très grande majorité les femmes et les filles. Une étude de 2019 a révélé que 96 % des vidéos deepfakes sont pornographiques et que tout le contenu des cinq sites Web les plus populaires pour les deepfakes pornographiques cible les femmes (282). La même étude a révélé que la grande majorité des abus deepfakes (99 % sur les sites pornographiques deepfakes et 81 % sur YouTube) visent les femmes artistes, suivies des femmes politiques (12 % sur YouTube). De plus, les deepfakes sexuels sont de plus en plus utilisés comme outil de violence conjugale, affectant de manière disproportionnée les femmes (271, 283). Une enquête représentative à l'échelle nationale auprès de 1 403 adultes britanniques a indiqué que les femmes étaient beaucoup plus susceptibles que les hommes de déclarer avoir peur de devenir une cible de pornographie deepfake, de devenir la cible d'une arnaque deepfake et de devenir la cible d'autres deepfakes potentiellement dangereux (284*). Cette préoccupation accrue des femmes pourrait refléter une prise de conscience de leur vulnérabilité accrue à de tels abus, suggérant un impact psychologique potentiel de cette technologie même sur celles qui ne sont pas directement ciblées. Cependant, la taille de l'échantillon de l'enquête était limitée et n'était pas représentative à l'échelle mondiale. De manière générale, des recherches supplémentaires sont nécessaires pour comprendre l'impact psychologique des deepfakes sur les femmes.

Les enfants sont confrontés à différents types de préjudices causés par le contenu sexuel généré par l'IA. Tout d'abord, les acteurs malveillants peuvent exploiter les outils d'IA pour générer du CSAM. Fin 2023, une enquête universitaire a découvert des centaines d'images d'abus sexuels sur mineurs dans un ensemble de données ouvert utilisé pour former des modèles populaires de génération de texte en image par l'IA tels que Stable Diffusion (285). Au Royaume-Uni, parmi les adultes interrogés qui ont déclaré avoir été exposés à des deepfakes sexuels au cours des six derniers mois, 17 % pensaient avoir vu des images représentant des mineurs (286). Deuxièmement, les enfants peuvent aussi commettre des abus en utilisant l'IA. Au cours de l'année dernière, les écoles ont commencé à se débattre avec un nouveau problème : les élèves utilisent des « applications de nudification » facilement téléchargeables pour générer et distribuer des images pornographiques dénudées de leurs pairs (en majorité des femmes) (287).

Depuis la publication du rapport intermédiaire, de nouvelles données ont suggéré une prévalence significative du contenu généré par l'IA en ligne. Au Royaume-Uni, une étude a révélé que 43 % des personnes âgées de 16 ans et plus déclarent avoir vu au moins un deepfake (sous forme de vidéos, d'imitations vocales et d'images) en ligne au cours des six derniers mois (50 % chez les enfants âgés de 8 à 15 ans) (286). Cependant, les données fiables restent relativement limitées. Comprendre l'impact des deepfakes sur les individus nécessitera des recherches plus approfondies sur une période prolongée.

Les contre-mesures qui aident les gens à détecter les faux contenus générés par l'IA, comme les étiquettes d'avertissement et le filigrane, montrent une efficacité mitigée. Certains outils d'IA peuvent aider à détecter les anomalies dans les images et les signaler comme étant probablement du faux contenu généré par l'IA. Cela se fait soit en utilisant des algorithmes d'apprentissage automatique pour rechercher des caractéristiques spécifiques dans les fausses images, soit en entraînant des réseaux neuronaux profonds à identifier et analyser les caractéristiques anormales des images de manière indépendante (288). Les étiquettes d'avertissement sur les contenus potentiellement trompeurs ont montré une efficacité limitée, même dans des contextes moins dangereux - par exemple, dans une étude expérimentale utilisant des vidéos générées par l'IA d'une personnalité publique aux côtés de clips authentiques, les étiquettes d'avertissement n'ont amélioré le taux de détection des participants que de 10,7 % à 21,6 % (289). Cependant, l'écrasante majorité des répondants qui ont reçu des avertissements n'étaient toujours pas en mesure de distinguer les deepfakes des vidéos non modifiées (289). Une autre mesure d'authentification destinée à empêcher les faux contenus générés par l'IA est le « filigrane », qui consiste à intégrer une signature numérique dans le contenu lors de sa création. Les techniques de tatouage numérique se sont révélées prometteuses pour aider les gens à identifier l'origine et l'authenticité des médias numériques pour les vidéos (290, 291), les images (292, 293, 294*), l'audio (295, 296) et le texte (297). Cependant, les techniques de tatouage numérique sont confrontées à plusieurs limites, notamment la suppression du tatouage numérique par des adversaires sophistiqués (298*, 299) et les méthodes permettant de tromper les détecteurs de tatouage numérique (299). Il existe également des préoccupations concernant la confidentialité et l'utilisation abusive potentielle de la technologie de tatouage numérique pour suivre et identifier les utilisateurs (300). De plus, pour de nombreux types de contenus préjudiciables abordés dans cette section, tels que les contenus pornographiques ou intimes non consentis, la capacité d'identifier le contenu comme étant généré par l'IA n'empêche pas nécessairement le préjudice de se produire. Même lorsque le contenu s'avère être faux, les dommages à la réputation et aux relations peuvent persister, car les personnes conservent souvent leur réaction émotionnelle initiale au contenu – par exemple, la position d'un individu dans sa communauté peut ne pas être restaurée simplement en révélant que le contenu est faux.

Les décideurs politiques qui s'efforcent de limiter les dommages causés aux individus par les faux contenus générés par l'IA doivent relever plusieurs défis majeurs. Il est difficile d'évaluer l'ampleur du problème en raison du sous-déclaration et du manque de statistiques fiables. Cela peut compliquer la détermination de l'intervention appropriée. Les méthodes de détection et les techniques de tatouage actuelles, bien que progressant, affichent des résultats mitigés et se heurtent à des défis techniques persistants. Cela signifie qu'il n'existe actuellement aucune solution unique et robuste pour détecter et réduire la diffusion de contenus nuisibles générés par l'IA. Enfin, les progrès rapides de la technologie de l'IA dépassent souvent les méthodes de détection, ce qui met en évidence les limites potentielles du recours exclusif aux interventions techniques et réactives.

Pour les pratiques de gestion des risques liées au faux contenu généré par l'IA, voir :

- [3.4.1. Former des modèles plus fiables](#)
- [3.4.2. Suivi et intervention](#)

2.1.2. Manipulation de l'opinion publique

INFORMATIONS CLÉS

- Les acteurs malveillants peuvent utiliser l'IA à usage général pour générer du faux contenu tel que du texte, Les images ou vidéos sont utilisées pour tenter de manipuler l'opinion publique. Les chercheurs estiment que si elles réussissent, de telles tentatives pourraient avoir plusieurs conséquences néfastes.
- L'IA polyvalente peut générer du contenu potentiellement persuasif à une échelle sans précédent et avec un haut degré de sophistication. Auparavant, la génération de contenu pour manipuler l'opinion publique impliquait souvent un compromis important entre qualité et quantité. Cependant, les résultats de l'IA à usage général sont souvent impossibles à distinguer du contenu généré par des humains, et leur génération est extrêmement peu coûteuse. Certaines études ont également montré qu'ils étaient aussi convaincants que le contenu généré par des humains.
- Cependant, il n'existe pas de consensus scientifique sur l'impact attendu de cette utilisation abusive de l'IA à usage général. Il existe peu de preuves sur les effets sociétaux plus larges des fausses informations, qu'elles soient créées intentionnellement ou partagées à l'insu, et qu'elles soient ou non activées par l'IA. Certains chercheurs pensent que les tentatives de manipulation de l'opinion publique à l'aide de l'IA à usage général sont principalement entravées par le manque de canaux de distribution efficaces. Ce point de vue implique que les progrès dans la génération de contenu manipulateur devraient avoir un impact limité sur l'efficacité de ces campagnes.
- Depuis la publication du rapport intermédiaire (mai 2024), de nouvelles recherches ont été menées sur la viralité des tentatives de manipulation basées sur l'IA et sur les mesures d'atténuation possibles. Une nouvelle étude révèle que le contenu manipulateur généré par l'IA est perçu comme moins précis mais partagé à des taux similaires au contenu généré par l'homme, ce qui suggère que ce contenu peut facilement devenir viral, qu'il soit généré par l'IA ou par l'homme. De nouvelles méthodes de détection techniques intégrant à la fois des données textuelles et visuelles ont montré un certain succès, mais ne sont pas entièrement fiables.
- Les décideurs politiques sont confrontés à des techniques d'atténuation limitées et à des compromis difficiles. Dans certains contextes, il peut être difficile de concilier le risque de manipulation des adresses par l'IA à usage général avec la protection de la liberté d'expression. En outre, à mesure que les résultats de l'IA à usage général deviennent de plus en plus convaincants et réalistes, la détection des cas de manipulation par l'IA peut devenir plus difficile. Les techniques de prévention, telles que le tatouage numérique du contenu, sont utiles mais peuvent être contournées avec un effort modéré.

Définitions clés

- Contenu factice généré par l'IA : contenu audio, texte ou visuel, produit par l'IA générative, qui représente des personnes ou des événements d'une manière qui diffère de la réalité, de manière malveillante ou trompeuse, par exemple en montrant des personnes faisant des choses qu'elles n'ont pas faites, en disant des choses qu'elles n'ont pas dites, en changeant le lieu d'événements réels ou en représentant des événements qui ne se sont pas produits.

- Agent IA : une IA à usage général qui peut élaborer des plans pour atteindre des objectifs, effectuer de manière adaptative des tâches impliquant plusieurs étapes et des résultats incertains en cours de route et interagir avec son environnement (par exemple en créant des fichiers, en effectuant des actions sur le Web ou en déléguant des tâches à d'autres agents) avec peu ou pas de surveillance humaine.

L'IA à usage général peut aider les gens à générer du contenu réaliste à grande échelle, que des acteurs malveillants pourraient utiliser pour tenter de manipuler l'opinion publique et diffuser certains récits. Des études montrent que les humains trouvent souvent impossible de distinguer le texte généré par l'IA à usage général du contenu authentique généré par l'homme (301, 302, 303, 304). En outre, les recherches indiquent que même si les gens ont du mal à identifier avec précision le contenu généré par l'IA, ils surestiment souvent leur capacité à le faire (305). Il existe également des preuves que ce type de contenu est déjà diffusé à grande échelle (306).

Des recherches récentes ont observé une augmentation significative des articles de presse générés par l'IA (307) et ont découvert que les modèles linguistiques de l'IA peuvent réduire les coûts de génération de contenu jusqu'à 70 % pour les modèles hautement fiables (308*).

Il existe des preuves que le contenu généré par l'IA à usage général peut être aussi persuasif que le contenu généré par des humains, du moins dans des conditions expérimentales. Des travaux récents ont mesuré la force de persuasion des messages politiques générés par l'IA à usage général. Plusieurs études ont montré qu'ils peuvent influencer l'opinion des lecteurs d'expériences psychologiques (309, 310, 311, 312, 313*), de manière potentiellement durable (314), bien que la généralisabilité de ces effets aux contextes du monde réel reste sous-étudiée. Une étude a révélé que lors des débats, les gens étaient tout aussi susceptibles d'être d'accord avec les opposants à l'IA qu'avec leurs opposants humains (315), et plus susceptibles d'être persuadés par l'IA si celle-ci avait accès à des informations personnelles du type de celles que l'on peut trouver sur les comptes de médias sociaux.

Des recherches récentes explorent également la manière dont les agents d'IA à usage général pourraient influencer les croyances des utilisateurs à l'aide de techniques plus sophistiquées, notamment en créant et en exploitant la dépendance émotionnelle des utilisateurs, en alimentant leurs angoisses ou leur colère, ou en menaçant de divulguer des informations si les utilisateurs ne se conforment pas (316*).

À mesure que les systèmes d'IA à usage général gagnent en capacité, il est prouvé qu'il deviendra plus facile de les utiliser de manière malveillante à des fins trompeuses ou manipulatoires, peut-être même avec une efficacité supérieure à celle des humains qualifiés, et d'encourager les utilisateurs à prendre des mesures qui vont à l'encontre de leurs propres intérêts (317, 318*). Il existe également des preuves que les systèmes d'IA peuvent utiliser de nouvelles tactiques de manipulation spécifiques à l'IA auxquelles les humains sont particulièrement vulnérables, car nos défenses contre la manipulation ont été développées en réponse à d'autres humains, et non à des IA (319). Cependant, les systèmes d'IA peuvent également jouer un rôle déterminant dans l'atténuation de la persuasion alimentée par l'IA.

Cependant, il existe un débat général concernant l'impact des tentatives de manipulation de l'opinion publique, que ce soit à l'aide d'une IA à usage général ou non. Une revue systématique des études empiriques pertinentes sur les fausses nouvelles a révélé que seules huit des 99 études examinées tentaient de mesurer les impacts directs (320). Ces études ont généralement constaté que la diffusion et la consommation de fausses nouvelles étaient limitées et fortement concentrées sur des groupes d'utilisateurs spécifiques, ce qui jette le doute sur les hypothèses antérieures concernant son influence généralisée sur les résultats des élections. Cependant, ces résultats n'indiquent pas nécessairement une grande résilience aux tentatives de manipulation et de persuasion, et les fausses nouvelles peuvent avoir une portée plus large ou plus large.

Les effets indésirables de l'IA vont au-delà de son objectif initial. Certaines études suggèrent que si les gens peuvent théoriquement distinguer les informations vraies des informations fausses, ils n'ont souvent pas la motivation pour le faire, se concentrant plutôt sur leurs motivations personnelles ou sur la maximisation de l'engagement sur les médias sociaux (321, 322, 323). Indépendamment du débat universitaire sur l'efficacité, l'inquiétude du public face aux tentatives de manipulation de l'opinion publique menées par l'IA reste élevée – par exemple, une enquête de 2024 a révélé qu'une majorité d'Américains de tous les horizons politiques étaient très préoccupés par le fait que l'IA soit utilisée pour créer de fausses informations sur les candidats aux élections (324). Cependant, ce résultat n'est peut-être pas représentatif des attitudes mondiales.

En outre, il n'existe pas de consensus sur la question de savoir si la génération de faux contenus plus réalistes à grande échelle doit conduire à des campagnes de manipulation plus efficaces, ou si le principal obstacle à de telles campagnes est la distribution (voir la figure 2.1). Certains experts ont fait valoir que le principal obstacle pour les acteurs qui tentent d'avoir un impact à grande échelle avec du faux contenu n'est pas la génération du contenu, mais sa distribution à grande échelle (325). De même, certaines recherches suggèrent que les « cheapfakes » (méthodes moins sophistiquées de manipulation de contenu audiovisuel qui ne dépendent pas de l'utilisation d'IA à usage général) pourraient être aussi nuisibles que les deepfakes plus sophistiqués (326). Si cela est vrai, cela étayerait l'hypothèse selon laquelle la qualité du faux contenu est actuellement moins décisive pour le succès d'une campagne de manipulation à grande échelle que les défis liés à la distribution de ce contenu à de nombreux utilisateurs. Les plateformes de médias sociaux peuvent utiliser diverses techniques pour réduire la portée du contenu susceptible d'être de cette nature. Ces techniques sont souvent relativement efficaces, mais leur impact sur la liberté d'expression suscite des inquiétudes. Elles comprennent la modération humaine du contenu, l'étiquetage du contenu potentiellement trompeur et l'évaluation de la crédibilité de la source. Dans le même temps, les recherches montrent depuis des années que les algorithmes des médias sociaux privilégient souvent l'engagement et la viralité plutôt que l'exactitude ou l'authenticité du contenu, ce qui, selon certains chercheurs, pourrait favoriser la diffusion rapide du contenu généré par l'IA pour manipuler l'opinion publique (327).

Les chercheurs ont également exprimé des inquiétudes plus larges quant à l'érosion de la confiance dans l'environnement informationnel à mesure que le contenu généré par l'IA devient plus répandu. Certains chercheurs craignent qu'à mesure que les capacités d'IA à usage général se développent et soient de plus en plus utilisées pour générer et diffuser des messages à grande échelle, qu'ils soient exacts, intentionnellement faux ou non, les gens puissent en venir à se méfier davantage de l'information, ce qui pourrait poser de graves problèmes pour la délibération publique. Les acteurs malveillants pourraient exploiter cette perte de confiance généralisée en niant la véracité de preuves réelles et défavorables, en prétendant qu'elles sont générées par l'IA – un phénomène connu sous le nom de « dividende des menteurs » (328, 329). Cependant, la société pourrait également s'adapter rapidement aux changements induits par l'IA dans l'environnement de l'information. Dans ce scénario plus optimiste, les individus pourraient adapter leurs normes communes pour déterminer si une information ou une source est crédible ou non. La société s'est ainsi adaptée aux changements technologiques passés, comme l'introduction de logiciels de retouche d'images traditionnelles.

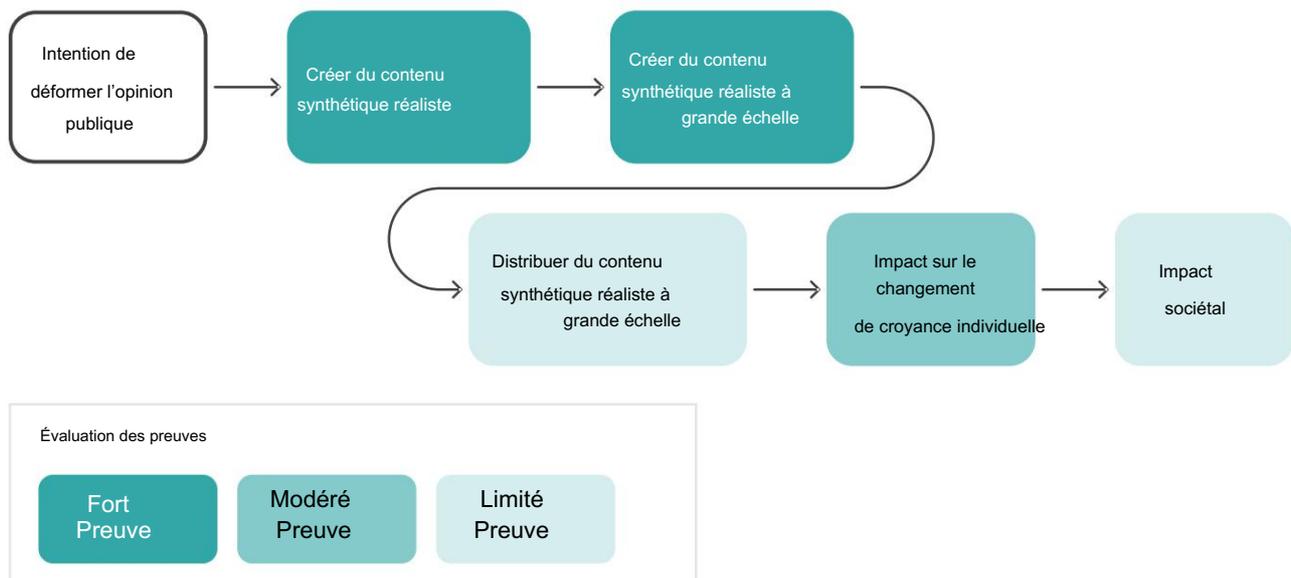


Figure 2.1 : Plusieurs étapes séparent l'intention initiale de manipuler l'opinion publique de l'impact potentiel sur la société. Bien qu'il existe de solides preuves de la capacité technique de créer du contenu généré par l'IA, les preuves deviennent rares aux stades ultérieurs, reflétant des lacunes de recherche plutôt qu'une inefficacité prouvée de telles campagnes. Il convient de noter que les impacts sociétaux peuvent également se produire par d'autres mécanismes que ceux décrits ici, comme une érosion générale de la confiance dans les sources d'information, même sans changements mesurables dans les croyances individuelles. Source : International AI Safety Report.

Depuis la publication du rapport intermédiaire, de nouvelles perspectives ont émergé concernant le contenu généré par l'IA. Une étude expérimentale récente a révélé que si les gens perçoivent les fausses nouvelles générées par l'IA comme moins précises que celles générées par l'homme (d'environ 20 %), ils partagent les deux types de fausses nouvelles dans des proportions similaires (environ 12 %), soulignant que le contenu fabriqué, qu'il soit généré par l'IA ou par l'homme, peut facilement devenir viral (330). Dans l'expérience, près de 99 % des sujets de l'étude n'ont pas réussi à identifier les fausses nouvelles générées par l'IA au moins une fois, ce que les auteurs ont attribué à la capacité des grands LLM de pointe à imiter le style et le contenu de sources fiables. De nouvelles méthodes de détection ont réussi à combiner l'analyse textuelle et visuelle, remédiant aux limites antérieures des approches n'utilisant qu'un seul type de données, comme uniquement du texte ou uniquement des images (331).

Les techniques actuelles d'identification du contenu généré par l'IA à usage général sont utiles mais souvent faciles à contourner. Les chercheurs ont utilisé diverses méthodes pour identifier la paternité potentielle de l'IA (332, 333). Les techniques d'« analyse de contenu » explorent les propriétés statistiques du texte, telles que les fréquences inhabituelles de caractères ou les distributions de longueur de phrase incohérentes, qui s'écartent des modèles généralement observés dans l'écriture humaine (334, 335, 336). Les techniques d'« analyse linguistique » examinent les éléments stylistiques, tels que la reconnaissance des sentiments ou des entités nommées, pour découvrir des incohérences ou des modèles de langage non naturels indiquant une génération d'IA (337, 338). Les chercheurs peuvent parfois également détecter le texte généré par l'IA en mesurant sa lisibilité, car l'écriture de l'IA présente souvent des modèles inhabituels par rapport à l'écriture humaine (339). Cependant, tous les contenus générés par l'IA ne sont pas des fausses nouvelles, et certaines recherches révèlent un biais intéressant dans les outils de détection de fausses nouvelles : ils ont tendance à classer de manière disproportionnée le contenu généré par LLM comme des fausses nouvelles, même lorsqu'il est véridique (340). Une étude portant sur sept détecteurs de contenu d'IA largement utilisés a identifié une autre limitation potentielle de ces derniers

outils : ils ont affiché un parti pris contre les auteurs non anglophones, classant souvent à tort leur travail comme généré par l'IA (341). Enfin, les chercheurs en IA ont également proposé d'autres approches pour détecter le contenu généré par l'IA, comme le « tatouage », dans lequel une signature invisible identifie le contenu numérique comme généré ou modifié par l'IA. Le tatouage peut aider à la détection du contenu généré par l'IA, mais peut généralement être contourné par des acteurs moyennement sophistiqués, comme cela est expliqué dans la [section 2.1.1. Dommages causés aux individus par le biais de faux contenus](#).

Les premières expériences démontrent que la collaboration humaine avec l'IA peut améliorer la détection de textes générés par l'IA. Cette approche a augmenté la précision de détection de 6,36 % pour les non-experts et de 12,76 % pour les experts par rapport aux efforts individuels dans une étude récente (342). Bien que la détection collaborative purement humaine ne soit probablement pas évolutive pour traiter la grande quantité de contenu généré quotidiennement, la recherche reste précieuse. Par exemple, les données issues de la collaboration humaine peuvent être utilisées pour former et améliorer les systèmes de détection de l'IA. De plus, pour les contenus particulièrement difficiles ou à enjeux élevés, la collaboration humaine peut compléter la détection de l'IA. Cependant, l'effet à long terme de ces efforts collaboratifs sur la résilience du public aux tentatives de manipulation reste à déterminer, et d'autres études sont nécessaires pour valider ces premiers résultats.

Les décideurs politiques qui s'efforcent de réduire le risque de manipulation de l'opinion publique par l'IA doivent relever plusieurs défis, notamment en tentant d'atténuer ces risques en protégeant la liberté d'expression (343, 344) et en déterminant les cadres de responsabilité juridique appropriés (345, 346, 347).

Les décideurs politiques sont également confrontés à une incertitude quant à l'impact réel des campagnes de manipulation, compte tenu des preuves mitigées sur leur efficacité et des données limitées sur leur prévalence (voir figure 2.1). Un autre défi est l'évolution continue de l'IA, les comportements adaptatifs des utilisateurs et les améliorations continues des systèmes d'IA, qui créent un cycle perpétuel d'adaptation et de contre-adaptation entre les méthodes de détection et le contenu généré par l'IA.

Pour les pratiques de gestion des risques liées à la manipulation de l'opinion publique, voir :

- [3.3. Identification et évaluation des risques](#)
- [3.4.2. Suivi et intervention](#)

2.1.3. Cyberinfraction

INFORMATIONS CLÉS†

- Les attaquants commencent à utiliser l'IA à usage général pour des opérations cybernétiques offensives, Les risques sont croissants mais limités pour l'instant. Les systèmes actuels ont démontré leurs capacités dans des tâches de cybersécurité de complexité faible et moyenne, les acteurs malveillants parrainés par l'État explorant activement l'IA pour surveiller les systèmes cibles. Des acteurs malveillants de différents niveaux de compétence peuvent exploiter ces capacités contre des personnes, des organisations et des infrastructures critiques telles que les réseaux électriques.
- Le risque cybernétique est dû au fait que l'IA polyvalente permet des opérations rapides et parallèles à grande échelle et réduit les obstacles techniques. Si les connaissances spécialisées restent essentielles, les outils d'IA réduisent l'effort et les connaissances humaines nécessaires pour surveiller les systèmes cibles et obtenir un accès non autorisé.
- L'IA polyvalente offre d'importantes capacités cybernétiques à double usage. Les données montrent qu'elle pourrait accélérer des processus tels que la découverte de vulnérabilités, qui sont essentielles pour lancer des attaques ainsi que pour renforcer les défenses. Cependant, les contraintes de ressources et les réglementations peuvent empêcher les services critiques et les petites organisations d'adopter des défenses renforcées par l'IA. L'impact final de l'IA sur l'équilibre entre attaquant et défenseur reste incertain.
- Depuis la publication du rapport intermédiaire (mai 2024), les systèmes d'IA à usage général ont montré des progrès significatifs dans l'identification et l'exploitation des vulnérabilités informatiques. Les systèmes d'IA ont trouvé et exploité de manière autonome des vulnérabilités dans de véritables projets de logiciels open source. Des prototypes de recherche récents ont trouvé et exploité de manière autonome des vulnérabilités qui prennent quelques minutes aux équipes de sécurité humaines les plus rapides à trouver, mais qui ont du mal à gérer des scénarios plus complexes. L'IA à usage général a également été utilisée pour trouver et corriger une vulnérabilité exploitable jusqu'alors inconnue dans un logiciel largement utilisé (SQLite).
- En principe, le risque semble au moins partiellement gérable, mais il reste des défis majeurs à relever en matière d'évaluation. Les progrès rapides des capacités rendent difficile l'exclusion de risques à grande échelle à court terme, ce qui souligne la nécessité d'évaluer et de surveiller ces risques. De meilleures mesures sont nécessaires pour comprendre les scénarios d'attaque réels, en particulier lorsque les humains et l'IA travaillent ensemble. L'un des défis majeurs consiste à atténuer les capacités offensives sans compromettre les applications défensives.

Définitions clés

- Logiciel malveillant : logiciel nuisible conçu pour endommager, perturber ou obtenir un accès non autorisé à un système informatique. Il comprend les virus, les logiciels espions et autres programmes malveillants qui peuvent voler des données ou causer des dommages.

† Veuillez vous référer à la [mise à jour du Président](#) sur les dernières avancées en matière d'IA après la rédaction de ce rapport.

- Ransomware : un type de logiciel malveillant qui verrouille ou crypte les fichiers ou le système d'un utilisateur, les rendant inaccessibles jusqu'à ce qu'une rançon (généralement de l'argent) soit versée à l'attaquant.
- Piratage informatique : acte consistant à exploiter les vulnérabilités ou les faiblesses d'un système informatique, d'un réseau ou d'un logiciel pour obtenir un accès non autorisé, manipuler des fonctionnalités ou extraire des informations.
- Tests de pénétration : pratique de sécurité dans laquelle des experts agréés ou des systèmes d'IA simulent des cyberattaques sur un système informatique, un réseau ou une application afin d'évaluer proactivement sa sécurité. L'objectif est d'identifier et de corriger les faiblesses avant qu'elles ne soient exploitées par de véritables attaquants.
- Défis CTF (Capture the Flag) : exercices souvent utilisés dans les formations en cybersécurité, conçus pour tester et améliorer les compétences des participants en les mettant au défi de résoudre des problèmes liés à la cybersécurité, comme trouver des informations cachées ou contourner les défenses de sécurité.
- Vulnérabilité zero-day : une faille de sécurité non découverte ou non corrigée dans un logiciel ou un matériel. Comme les attaquants peuvent déjà l'exploiter, les développeurs ont « zéro jour » pour le corriger.
- Porte dérobée matérielle : une fonctionnalité d'un appareil, intentionnellement ou non, créée par un fabricant ou tiers, qui peuvent être utilisés pour contourner les protections de sécurité afin de surveiller, contrôler ou extraire des données à l'insu de l'utilisateur.

Les cyberopérations offensives impliquent généralement la conception et le déploiement de logiciels malveillants (malware) et l'exploitation des vulnérabilités des systèmes logiciels et matériels, ce qui conduit à de graves failles de sécurité. Une chaîne d'attaque standard commence par la reconnaissance du système cible, suivie d'une découverte itérative, de l'exploitation des vulnérabilités et de la collecte d'informations supplémentaires. Ces actions nécessitent une planification minutieuse et une exécution stratégique pour atteindre les objectifs de l'adversaire tout en évitant d'être détecté. Certains experts craignent que l'IA à usage général puisse améliorer ces opérations en automatisant la détection des vulnérabilités, en optimisant les stratégies d'attaque et en améliorant les techniques d'évasion (348, 349). Ces capacités avancées profiteraient à tous les attaquants. Par exemple, les acteurs étatiques pourraient les exploiter pour cibler les infrastructures nationales critiques (CNI), ce qui entraînerait des perturbations généralisées et des dommages importants. Dans le même temps, l'IA à usage général pourrait également être utilisée à des fins défensives, par exemple pour trouver et corriger les vulnérabilités.

L'IA à usage général peut aider à recueillir des informations, réduisant ainsi l'effort humain. Par exemple, dans les attaques par rançongiciel, les acteurs malveillants effectuent d'abord manuellement une reconnaissance offensive et exploitent les vulnérabilités pour pénétrer dans le réseau cible, puis diffusent des logiciels malveillants qui se propagent sans intervention humaine (350). La phase d'entrée est souvent techniquement difficile et sujette à l'échec. L'IA à usage général est explorée par des attaquants parrainés par l'État comme une aide pour accélérer le processus (351*, 352*). Cependant, bien qu'il existe des systèmes à usage général qui ont effectué la découverte de vulnérabilités de manière autonome (voir les paragraphes suivants), les systèmes publiés n'ont pas encore exécuté de manière autonome des intrusions réelles dans des réseaux et des systèmes – des tâches qui sont par nature plus complexes.

Les systèmes d'IA à usage général se sont considérablement améliorés dans la détection autonome des vulnérabilités cybernétiques

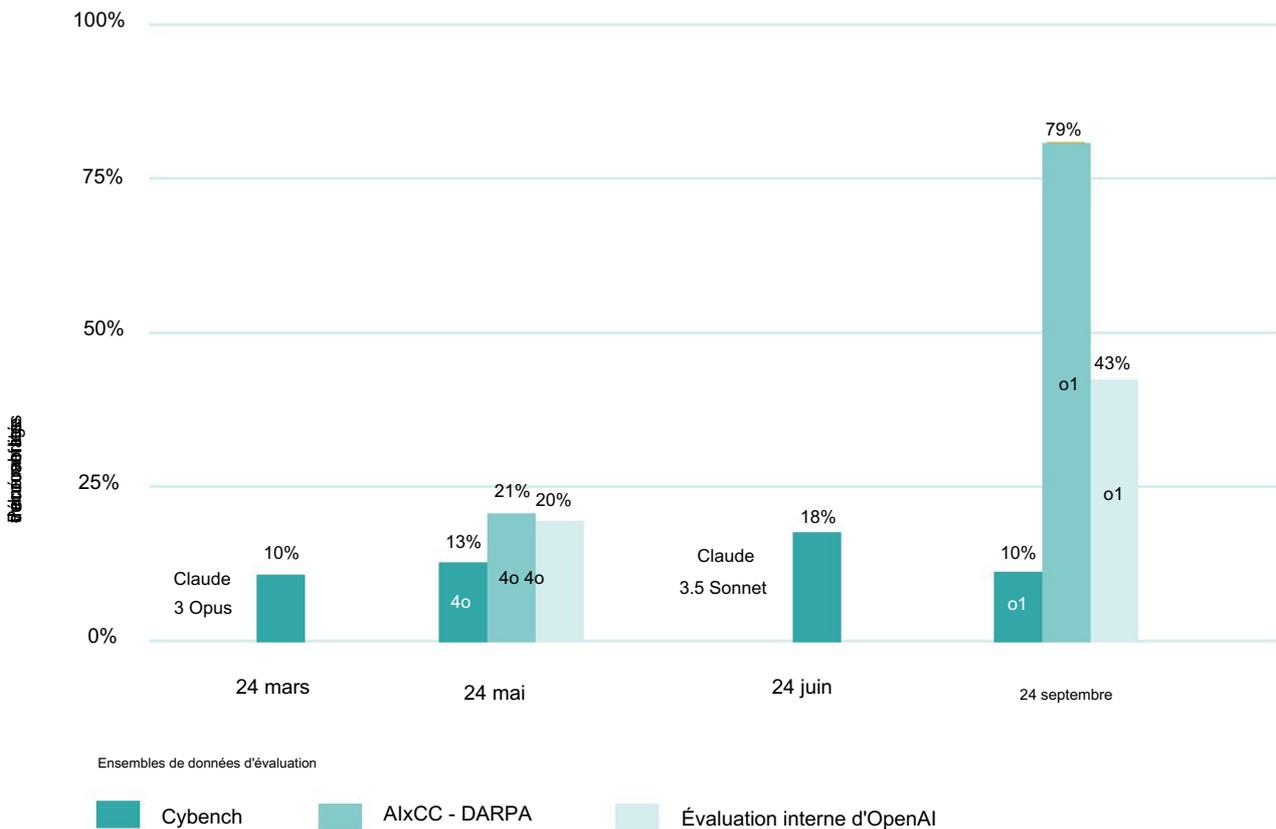


Figure 2.2 : Les progrès récents dans la capacité des modèles d'IA à trouver et exploiter les vulnérabilités de cybersécurité de manière autonome se sont accrus dans plusieurs tests de référence. Dans le cadre du AI Cyber Challenge de la DARPA et de l'ARPA-H (353, 359), le nouveau modèle o1 d'OpenAI (septembre 2024) a largement surpassé GPT-4o (mai 2024), détectant de manière autonome 79 % des vulnérabilités contre 21 % pour GPT-4o. Les tests sur Cybench (358) ont montré que les taux de détection des vulnérabilités s'amélioraient de 10 % (Claude 3 Opus, mars 2024) à 17,5 % (Claude 3.5 Sonnet, juin 2024). Les évaluations internes du concours de piratage CTF d'OpenAI au niveau du secondaire sont passées de 20 % à 43 %, bien que les modèles aient toujours du mal à effectuer des tâches de référence plus complexes (2*). Seul le nouveau modèle le plus performant pour chaque mois est présenté. Sources : Agence des projets de recherche avancée de défense, 2024 (353) ; Ristea et al., 2024 (359) ; Zhang et al., 2024 (358) ; OpenAI, 2024 (2*).

L'IA à usage général peut aider les attaquants à découvrir les vulnérabilités dans le code source dans une certaine mesure, mais les méthodes traditionnelles restent dominantes pour l'instant. Dans cette tâche, l'analyste examine le code source d'un projet logiciel (comme un serveur Web open source ou un pare-feu) pour identifier les failles de sécurité exploitables.

Depuis la publication du rapport intermédiaire, les capacités cybernétiques de l'IA à usage général en matière de découverte de vulnérabilités ont considérablement augmenté. Lors du défi AlxCC de la DARPA (353) les participants ont développé des systèmes capables de trouver, d'exploiter et de corriger de manière autonome les vulnérabilités dans de véritables projets de logiciels open source à l'aide d'une IA à usage général (354, 355, 356). La figure 2.2 illustre l'amélioration significative des performances des modèles d'IA à usage général pour détecter et exploiter (et parfois corriger) les vulnérabilités informatiques. De plus, Big Sleep de Google a été utilisé pour découvrir une vulnérabilité exploitable jusqu'alors inconnue dans le

Le logiciel open source SQLite (357*) est largement utilisé. Dans ce cas, la découverte a été utilisée pour corriger plutôt que pour exploiter la vulnérabilité. Les tests de pénétration (métriques) ont également considérablement progressé, fournissant un signal beaucoup plus clair des capacités des modèles et de leur amélioration au fil du temps (358).

L'IA à usage général a montré un succès faible à modéré dans l'automatisation du piratage des systèmes et des réseaux. Contrairement à la détection automatisée des vulnérabilités dans les logiciels, où les systèmes d'IA s'appuient sur l'accès au code source, le piratage est plus difficile pour l'IA, car elle doit exécuter chaque étape d'une attaque avec peu ou pas de connaissance préalable du fonctionnement interne du système ciblé (par exemple, collecter des informations sur la cible, trouver des points d'entrée, pénétrer dans le système, se déplacer dans le système et atteindre son objectif). Les opérations de piratage dans le monde réel nécessitent des actions exploratoires et des ajustements itératifs pour comprendre le fonctionnement d'un système cible, impliquant souvent des tests d'hypothèses et des changements de stratégie de manière dynamique (360). Ces tâches ont résisté à l'automatisation complète car elles nécessitent un niveau de précision exceptionnel (un seul caractère incorrect dans une entrée peut entraîner l'échec de l'approche dans son ensemble) et impliquent la résolution de plusieurs sous-tâches complexes sans conseils ni retours explicites.

Depuis la publication du rapport intermédiaire, les capacités d'attaque de cybersécurité de l'IA à usage général se sont améliorées, mais les modèles d'IA ne peuvent toujours pas battre les experts humains et ils ont du mal à gérer des scénarios plus complexes. Les défis CTF, où l'attaquant doit identifier et exploiter des vulnérabilités pour accéder à des données ou des systèmes protégés, sont devenus une référence typique en matière de cybersécurité. Avant le rapport intermédiaire (mai 2024), l'IA à usage général pouvait mener des attaques simples (127, 361, 362*) mais pas des attaques sophistiquées. Depuis lors, d'autres recherches ont réussi à obtenir de meilleurs résultats avec les systèmes d'IA. Par exemple, des équipes d'agents LLM peuvent collaborer efficacement pour trouver des vulnérabilités jusqu'alors inconnues (« zero-day »), bien que pas très compliquées (363). De plus, l'accès à de meilleurs outils (364) et l'introduction de modules permettant un raisonnement étape par étape (365) ont permis aux modèles d'IA à usage général de résoudre des tâches à partir de défis CTF établis (difficulté facile et moyenne). Cependant, sans ces aides au raisonnement, Google rapporte que son dernier modèle, Gemini 1.5, ne montre aucun gain de performance sur ses benchmarks CTF par rapport aux versions précédentes, et qu'il ne montre des améliorations que dans les tâches de cybersécurité offensives simples (49*). OpenAI rapporte que même si son récent modèle o1 s'améliore par rapport aux scores de référence de GPT-4o, il est toujours classé comme « à faible risque » dans ce domaine et reste dans des limites d'utilisation abusives gérables (2*). Divers modèles et collaborations entre modèles atteignent des performances comparables à celles des humains qui disposent d'environ 35 minutes par tâche : Sonnet 3.5 de base, o1-preview advising o1-preview (les deux versions) et o1-mini advising GPT-4o (129). Cette dynamique collaborative, où les modèles conseillent et affinent les résultats des uns et des autres, est de plus en plus utile pour les tâches à plusieurs étapes intolérantes aux erreurs telles que la cyberinfraction (voir également la section 1.2. [Capacités actuelles](#)). Les modèles sans aide n'ont pas été en mesure de résoudre les défis CTF qui ont pris aux meilleures équipes d'experts humains plus de 11 minutes de travail (358). Comme prévu, les modèles plus récents (par exemple GPT-4o et o1-preview d'OpenAI) fonctionnent mieux mais ont toujours du mal à générer des informations que les experts mettent plus de temps à comprendre.

Les systèmes d'IA à usage général peuvent réduire les connaissances techniques et l'expertise nécessaires pour mener à bien les différentes étapes de la chaîne d'attaque. Dans une chaîne d'attaque classique, un attaquant peut commencer par une reconnaissance pour identifier les vulnérabilités potentielles, utiliser une campagne de phishing pour obtenir un accès initial, obtenir des privilèges au sein du système cible, se déplacer latéralement sur le réseau et, finalement, exfiltrer des données sensibles ou déployer des ransomwares. En automatisant ou en aidant certaines parties de la chaîne d'attaque, l'IA à usage général réduit le besoin d'intervention d'experts, abaissant ainsi la barrière à l'entrée pour des attaques plus sophistiquées. Cependant, si l'IA peut accélérer le processus d'examen des informations accessibles au public, cela ne se traduit pas automatiquement par une expertise avancée. Dans des domaines tels que l'exploitation des vulnérabilités, l'IA à usage général peut aider, mais les experts doivent toujours intégrer des connaissances spécifiques au domaine pour rendre ces systèmes d'IA efficaces (353, 366*), un besoin qui n'a pas changé depuis le rapport intermédiaire.

Des groupes de hackers parrainés par des États auraient utilisé l'IA à usage général pour soutenir leurs activités de piratage. Par exemple, ces groupes ont utilisé l'IA à usage général pour traduire des documents techniques, analyser des vulnérabilités divulguées publiquement, rechercher des protocoles publics (par exemple, les communications par satellite), aider à la rédaction de scripts, résoudre des erreurs et développer des techniques d'évasion de détection pour les logiciels malveillants et les intrusions (351*).

L'IA à usage général ne devrait faire pencher la balance en faveur des attaquants que dans certaines conditions : 1. si l'IA à usage général automatise les tâches nécessaires à l'attaque mais pas les défenses correspondantes ; ou 2. si les capacités d'IA à usage général de pointe sont accessibles aux adversaires mais pas également disponibles pour tous les défenseurs. En particulier, les petites et moyennes entreprises (PME) peuvent ne pas être en mesure de se permettre des solutions de défense polyvalentes améliorées par l'IA. Par exemple, les hôpitaux, limités par des ressources de sécurité limitées et par la complexité de réseaux hérités hétérogènes, peuvent être plus lents à adopter des défenses basées sur l'IA, laissant leurs données hautement sensibles plus exposées aux cyberattaques sophistiquées. De même, les systèmes CNI (tels que les sous-stations électriques) ont souvent des critères stricts et sont prudents dans l'adoption de nouvelles technologies, y compris les défenses basées sur l'IA, en raison de problèmes de sécurité et d'exigences de gouvernance et/ou réglementaires. En revanche, les adversaires ne sont pas liés par de telles contraintes et peuvent adopter des capacités d'IA avancées plus rapidement.

Même si la détection pilotée par l'IA détecte les vulnérabilités dans le nouveau code avant qu'il n'atteigne la production, un défi majeur demeure : le code source déjà utilisé et antérieur à ces capacités. Une grande partie de ce code hérité n'a pas été examinée par des outils d'IA avancés, ce qui laisse des vulnérabilités potentielles non détectées. La correction de ces vulnérabilités après leur découverte est un processus lent, en particulier dans les environnements de production où les changements nécessitent des tests rigoureux pour éviter de perturber les opérations.

Par exemple, la vulnérabilité Heartbleed a continué à exposer les systèmes pendant des semaines après la mise à disposition d'un correctif, car les administrateurs ont dû faire face à des retards dans sa mise en œuvre (367). Cette situation créera potentiellement une période de transition critique, au cours de laquelle les défenseurs devront gérer et corriger le code ancien et non vérifié, tandis que les attaquants, libérés de telles contraintes et potentiellement équipés d'une IA avancée, pourront exploiter ces vulnérabilités avec moins d'efforts (une asymétrie de capacité). Au cours de cette transition, la disparité dans l'adoption de l'IA – en particulier parmi les PME et les systèmes d'infrastructure critiques qui sont plus lents

intégrer de nouvelles technologies comme l'IA – pourrait amplifier le déséquilibre entre les attaquants et les défenseurs.

Les contreparties défensives de certaines tâches offensives sont considérablement plus complexes, créant une asymétrie dans l'efficacité de l'IA à usage général lorsqu'elle est utilisée par les attaquants par rapport aux défenseurs.

Par exemple, les attaquants utilisant une IA à usage général peuvent intégrer furtivement des menaces au niveau du matériel (368) d'une manière que les défenseurs ont par nature du mal à prévoir ou à détecter. Ainsi, les attaquants contrôlent le degré de dissimulation et de complexité des vulnérabilités, tandis que les défenseurs doivent anticiper et détecter ces menaces délibérément masquées.

Le malware Stuxnet (369) a démontré comment de telles attaques peuvent causer des dommages physiques en ciblant les systèmes de contrôle industriels – il a perturbé les installations nucléaires iraniennes en manipulant les opérations matérielles. Bien qu'il n'existe aucune preuve publique que l'IA ait été utilisée pour automatiser et intensifier de telles menaces dans les systèmes de production, son impact potentiel sur la cybersécurité justifie une surveillance attentive. D'un autre côté, certaines applications d'IA pourraient également offrir des avantages asymétriques aux défenseurs. Par exemple, l'IA pourrait améliorer la sécurité des puces –

comme ceux utilisés dans les smartphones – en détectant et en atténuant les vulnérabilités pendant le processus de conception (370).

De plus, l'IA à usage général a déjà été intégrée dans les outils d'audit et de débogage (371*, 372*).

Les principales lacunes en matière de données probantes concernant les capacités actuelles de l'IA dans le domaine de la cybersécurité sont les suivantes :

- Évaluation complète des capacités : davantage d'études empiriques sont nécessaires pour évaluer les performances de l'IA dans des chaînes d'attaque complexes et réelles et pour suivre les tendances en matière de capacités, en particulier pour l'automatisation des attaques en plusieurs étapes. Les benchmarks existants, tels que les défis CTF, offrent des informations partielles, mais ne parviennent souvent pas à saisir l'étendue complète des capacités offensives pilotées par l'IA. Par exemple, l'évaluation comparative dans des environnements spécialisés, tels que les bancs d'essai d'infrastructures cyberphysiques, permettrait une évaluation plus réaliste de l'impact de l'IA dans des scénarios à enjeux élevés. En outre, l'absence de références de performance humaine rend difficile la contextualisation de la complexité des tâches en termes d'heures humaines, ce qui empêche de comparer avec précision les capacités de l'IA et des humains.
- Évaluation des attaques collaboratives homme-IA : la recherche sur la manière dont les attaquants pourraient exploiter l'IA aux côtés de l'expertise humaine est essentielle pour comprendre les avancées offensives potentielles. Des études devraient être menées pour déterminer comment l'IA peut améliorer les opérations menées par des humains dans des domaines tels que la prise de décision stratégique, l'allocation des ressources et les ajustements en temps réel, augmentant ainsi potentiellement l'efficacité et la sophistication des cyberattaques. De plus, les modèles d'IA produisent souvent des « quasi-accidents » que des humains ayant une expérience modérée de la cybersécurité pourraient facilement résoudre, ce qui suggère un avantage synergétique lorsque les humains et l'IA collaborent dans des efforts offensifs.

Les décideurs politiques qui se concentrent sur les cyber-risques seront confrontés à des défis, notamment celui d'évaluer de manière fiable les risques et les capacités de l'IA dans des contextes offensifs et défensifs. Les analyses comparatives des cyber-risques peuvent parfois surestimer les performances par rapport aux scénarios réels, car elles utilisent souvent des défis et du code provenant de plateformes telles que GitHub, que les modèles ont pu rencontrer lors de leur formation. Par conséquent, ces modèles peuvent déjà être familiarisés avec le code

ou ont bénéficié de tutoriels et de manuels de solutions trouvés dans des blogs et d'autres ressources en ligne. Cependant, les évaluations de capacités peuvent également être sous-estimées car il est difficile d'obtenir toutes les capacités d'un système ([1.2. Capacités actuelles](#)). De plus, les taux de réussite rapportés dans les tests de performance excluent généralement les quasi-échecs (cas où le modèle d'IA réussit presque l'attaque) (358), qui pourraient facilement être exploités par un opérateur humain pour mener à bien l'attaque.

Les décideurs politiques devront également relever des défis importants pour réglementer la recherche offensive en IA tout en préservant les capacités défensives. La recherche offensive en cybernétique est importante pour maintenir des défenses solides, et la restreindre pourrait affaiblir les stratégies de sécurité nationale, surtout si d'autres pays n'imposent pas de limitations similaires. Les décideurs politiques doivent peser les risques d'utilisation abusive par rapport aux avantages de cette recherche et trouver des moyens de réduire les risques d'utilisation abusive tout en protégeant les applications défensives (voir [3.3. Identification et évaluation des risques pour une discussion plus approfondie sur l'évaluation des risques et des capacités nuisibles](#)). Un autre problème crucial est la gestion des compromis impliqués dans la publication ouverte des pondérations des modèles d'IA à usage général, qui comporte à la fois des avantages importants et des risques d'utilisation abusive, comme exploré dans [2.4. Impact des modèles d'IA à usage général à pondération ouverte sur les risques liés à l'IA](#).

Pour les pratiques de gestion des risques liées aux cyberinfractions, voir :

- [3.3. Identification et évaluation des risques](#)
- [3.4.1. Former des modèles plus fiables](#)
- [3.4.2. Suivi et intervention](#)
- [3.4.3. Méthodes techniques de protection de la vie privée](#)

2.1.4. Attaques biologiques et chimiques

INFORMATIONS CLÉS†

- De plus en plus de preuves montrent que les avancées de l'IA à usage général sont bénéfiques pour la science tout en abaissant certains obstacles au développement d'armes chimiques et biologiques, tant pour les novices que pour les experts. De nouveaux modèles linguistiques peuvent générer des instructions techniques étape par étape pour créer des agents pathogènes et des toxines qui surpassent les plans rédigés par des experts titulaires d'un doctorat et font apparaître des informations que les experts ont du mal à trouver en ligne, bien que leur utilité pratique pour les novices reste incertaine. D'autres modèles démontrent des capacités à concevoir des protéines améliorées et à analyser quels agents pathogènes ou toxines candidats sont les plus nocifs. Les experts pourraient potentiellement les utiliser pour développer des armes plus avancées et des mesures défensives.
- L'impact réel de l'IA sur le développement et l'utilisation des armes, y compris en cas de pandémie
Les agents pathogènes restent flous en raison des exigences de confidentialité, des interdictions de tests et de la nécessité de meilleures évaluations. Les preuves clés concernant les acteurs malveillants, leurs goulots d'étranglement techniques et les évaluations de sécurité de l'IA relatives aux armes biologiques sont gardées confidentielles pour éviter toute utilisation abusive. Les tests sont souvent interdits en raison des graves dangers que représentent ces armes. D'autres évaluations sont nécessaires pour déterminer dans quelle mesure les systèmes actuels peuvent contribuer aux nombreuses étapes du développement des armes ; une expertise et des ressources substantielles restent des obstacles nécessaires.
- Ces derniers mois, les avancées ont permis de mieux cerner les risques et d'élargir les capacités biologiques de l'IA à usage général. Des efforts sont également en cours pour développer les meilleures pratiques en matière d'évaluation. Depuis le rapport intermédiaire (mai 2024), les modèles de langage à usage général ont fait des progrès substantiels dans les tests d'expertise en armes biologiques et de raisonnement scientifique général. L'IA a également démontré de nouvelles capacités dans la conception de protéines et dans le travail avec de multiples types de données scientifiques – notamment des produits chimiques, des protéines et de l'ADN – améliorant ainsi sa capacité à concevoir des structures biologiques complexes. Les implications en termes de risques sont toujours à l'étude, les premières preuves suggérant une augmentation des risques potentiels parallèlement aux avantages.
- Si les progrès rapides se poursuivent, cela créera des défis politiques urgents pour évaluer et gérer les risques biologiques. Les progrès rapides récents dans les critères de référence des risques rendent de plus en plus difficile d'exclure les risques à grande échelle dans les modèles à court terme. Les décideurs politiques doivent prendre des décisions avec des informations incomplètes et intégrer des recherches sur les menaces classifiées. À ces défis s'ajoutent les débats en cours sur les compromis risques-bénéfices de la publication de modèles à pondération ouverte, en particulier les outils d'IA pour créer des structures biologiques et chimiques, et le fait que les politiques qui dépendent des humains pour détecter les risques et intervenir pourraient être trop lentes pour s'adapter au rythme actuel de développement.

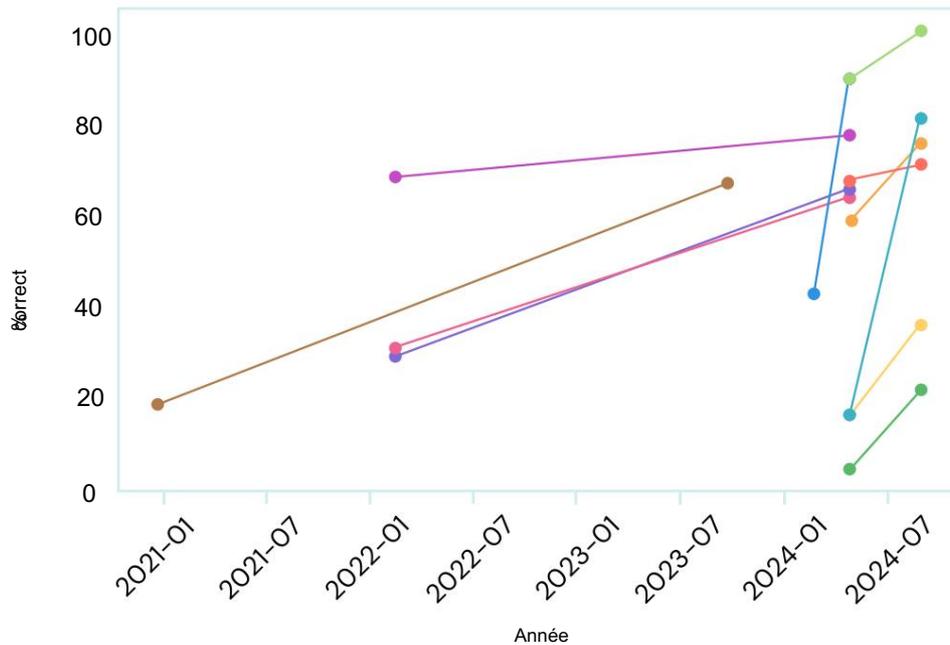
† Veuillez vous référer à la [mise à jour du Président](#) sur les dernières avancées en matière d'IA après la rédaction de ce rapport.

Définitions clés

- Science à double usage : recherche et technologie pouvant être appliquées à des fins bénéfiques, comme des médicaments ou des solutions environnementales, mais aussi potentiellement utilisées à mauvais escient pour causer du tort, comme dans le développement d'armes biologiques ou chimiques.
- Toxine : substance toxique produite par des organismes vivants (tels que des bactéries, des plantes ou des animaux), ou créée synthétiquement pour imiter une toxine naturelle, qui peut provoquer des maladies, des dommages ou la mort chez d'autres organismes en fonction de sa puissance et du niveau d'exposition.
- Pathogène : un micro-organisme, par exemple un virus, une bactérie ou un champignon, qui peut provoquer une maladie chez les humains, les animaux ou les plantes.
- Agent : Aux fins de la présente section, le terme « agent » désigne généralement un agent biologique, chimique ou substance toxicologique qui peut nuire aux organismes vivants. Les agents dans ce sens ne doivent pas être confondus avec les agents IA (voir ci-dessous).
- Agent IA : une IA à usage général qui peut élaborer des plans pour atteindre des objectifs, effectuer de manière adaptative des tâches impliquant plusieurs étapes et des résultats incertains en cours de route, et interagir avec son environnement (par exemple en créant des fichiers, en effectuant des actions sur le Web ou en déléguant des tâches à d'autres agents) avec peu ou pas de surveillance humaine.
- Biosécurité : un ensemble de politiques, de pratiques et de mesures (par exemple, diagnostics et vaccins) conçu pour protéger les humains, les animaux, les plantes et les écosystèmes contre les agents biologiques nocifs, qu'ils soient d'origine naturelle ou introduits intentionnellement.

Les risques associés à la science à double usage sont au cœur des préoccupations de la politique internationale de sécurité de l'IA. Cette section se concentre sur les armes chimiques et biologiques, mais il existe également des risques concernant les armes radiologiques et nucléaires. Ces armes de destruction massive, initialement développées dans le cadre de recherches scientifiques destinées à des fins pacifiques, illustrent le phénomène de la « science à double usage », où les innovations sont réutilisées pour des applications militaires. Parmi celles-ci, cette section se concentre sur les armes chimiques et biologiques, qui sont particulièrement préoccupantes en raison de la relative facilité d'obtention des matériaux nécessaires et de la large disponibilité des informations connexes. En conséquence, les risques liés aux armes biologiques ont occupé une place centrale dans les sommets sur la sécurité de l'IA et dans les discussions plus larges sur les impacts catastrophiques potentiels de l'IA avancée. En revanche, le risque que l'IA élargisse l'accès aux armes nucléaires et radiologiques est considéré comme plus faible, principalement en raison des obstacles importants à l'acquisition des matériaux requis. Cependant, l'implication de l'IA dans la prise de décision nucléaire introduirait des risques uniques. Certains experts craignent que la délégation du pouvoir de décision pour les lancements d'armes nucléaires à des systèmes d'IA n'augmente le risque d'erreurs critiques (voir [2.2.1. Problèmes de fiabilité](#)) ou de perte de contrôle (voir [2.2.3. Perte de contrôle](#)) (373). Les risques liés à la science à double usage s'étendent à d'autres avancées telles que les systèmes de navigation, la nanotechnologie, les robots et les drones autonomes, qui ont tous des applications militaires qui dépassent le cadre du présent rapport.

Les modèles d'IA sont récemment devenus plus performants dans les tâches à double usage et dans les tâches liées aux armes biologiques et chimiques.



Tâche de référence

Masters de LLM

- Création d'armes biologiques : résolution des problèmes (% de questions correctes)
- Création d'armes biologiques : tâches de laboratoire (% de questions correctes)
- Acquisition d'armes biologiques (% de questions correctes)
- Grossissement des armes biologiques (% de questions correctes)
- Formulation d'armes biologiques (% de questions correctes)
- Lancement d'armes biologiques (% de questions correctes)

Modèles biologiques à usage général

- Protéines se liant à de petites molécules (% de structures correctes)
- Protéines se liant à l'ADN (% structures correctes)
- Liaison des protéines aux protéines (% de structures correctes)
- Protéines se liant aux anticorps (% structures correctes)

Modèles biologiques spécialisés

- Prédire les mutations pandémiques courantes (% de mutations prédites)

Figure 2.3 : Les capacités à double usage en biologie ont augmenté au fil du temps pour les LLM (2*), l'IA biologique à usage général comme AlphaFold3 (23) et les modèles spécialisés (non à usage général) relatifs aux agents pathogènes (390). Ce graphique montre les scores de performance, calculés en pourcentage de précision pour les résultats récemment publiés par rapport aux résultats de pointe précédents. Les avancées récentes dans les LLM ont été particulièrement rapides, en comparant GPT4o (sortie en mai 2024) à o1 (sortie en septembre 2024). Français Les avancées notables sont la précision des LLM pour répondre aux questions sur la libération d'armes biologiques, qui est passée de 15 % à 80 %, et la capacité des IA biologiques à prédire comment les protéines interagissent avec les petites molécules (y compris dans les médicaments et les armes chimiques), qui est passée de 42 % à 90 % en 2024. En raison d'un manque de références standardisées et d'incohérences dans la manière dont la précision est calculée, les comparaisons sont limitées à quelques tâches et ne sont pas répétées de manière cohérente au fil du temps. Sources : OpenAI, 2024 (2*) (pour les LLM) ; Abramson et al., 2024 (23) (comparaison d'AlphaFold3 avec l'état de l'art précédent) ; Thadani et al., 2023 (390) (pour les modèles spécialisés relatifs aux agents pathogènes).

Certaines IA à usage général ont été développées spécifiquement pour les domaines scientifiques, offrant des capacités générales de compréhension et de conception de produits chimiques, d'ADN et de protéines. Les modèles formés à partir de données scientifiques ont des capacités variées, allant d'applications étroites telles que la prédiction de la structure des protéines, à une variété de capacités de prédiction et de conception. Dans ce rapport, les modèles largement capables formés à partir de données scientifiques sont inclus dans la définition de l'IA à usage général. Cependant,

Il existe un débat important au sein des communautés de l'IA et de la biologie concernant le point à partir duquel un modèle formé sur des données scientifiques peut être qualifié de « modèle à usage général » (voir [Introduction](#) pour une définition) ou de « modèle de base » (45). Par exemple, AlphaFold2 a été conçu pour la tâche étroite de prédiction de la structure des protéines, mais, grâce à un réglage fin, il s'est avéré applicable à une grande variété d'autres tâches, telles que la prédiction des interactions protéiques, la prédiction des petits sites de liaison moléculaire et la prédiction et la conception de peptides cycliques (374). Pour ces raisons, il satisfait à la définition d'un modèle d'IA à usage général de ce rapport. AlphaFold3 a été capable d'accomplir ces tâches avec une plus grande précision et sur une plus large gamme de molécules, même sans réglage fin (23). Ces outils d'IA à vocation scientifique amplifient le potentiel d'innovation chimique et biologique en accélérant la découverte scientifique, en optimisant la production et en permettant la conception précise de nouveaux éléments biologiques. Ils offrent également des opportunités prometteuses pour développer de nouveaux médicaments et mieux lutter contre les maladies infectieuses (375, 376). Ces outils ont généré des avancées scientifiques substantielles, suffisantes pour valoir à leurs créateurs le prix Nobel de chimie (377).

Le caractère à double usage des progrès scientifiques présente des risques complexes, car les innovations destinées à des fins bénéfiques, comme la médecine, ont historiquement conduit à la création d'armes chimiques et biologiques (378, 379). La grande majorité des dommages causés par les toxines et les maladies infectieuses résultent d'événements naturels, ce qui a donné lieu à des recherches approfondies pour aider à combattre ces menaces. Le développement et le déploiement intentionnels d'armes biologiques ont été éclairés par ces recherches, mais posent des difficultés considérables (380, 381). Beaucoup pensent que les progrès dans la conception, l'optimisation et la production de produits chimiques et biologiques, en partie dus à l'IA, ont peut-être facilité le développement d'armes chimiques et biologiques (382, 383, 384, 385). Les preuves présentées dans cette section suggèrent que l'IA à usage général amplifie les risques liés aux armes en aidant les novices (généralement définis comme des personnes titulaires d'une licence ou moins dans une discipline pertinente) à créer ou à accéder aux armes biologiques et chimiques existantes, et en permettant aux experts (généralement une personne titulaire d'un doctorat ou plus dans une discipline pertinente) de concevoir des armes plus dangereuses ou ciblées, ou de créer des armes existantes avec moins d'efforts.

Depuis la publication du rapport intermédiaire, la capacité des modèles d'IA à usage général à raisonner et à intégrer différents types de données s'est améliorée, et des progrès ont été réalisés dans la formulation de meilleures pratiques en matière de biosécurité. Plusieurs modèles ont été publiés depuis le rapport intermédiaire (mai 2024) qui intègrent différents types de données scientifiques ; un modèle de base pour les données scientifiques, AlphaFold 3, peut prédire la structure et les interactions entre une gamme de molécules, y compris les produits chimiques, l'ADN et les protéines, avec une plus grande précision que l'état de l'art précédent (voir la figure 2.3) (23), et un autre, ESM3, peut modéliser simultanément la séquence, la structure et la fonction des protéines (386*). Ces développements ouvrent de nouvelles possibilités pour la conception de produits biologiques qui ne ressemblent pas fortement aux produits naturels (387*). Le o1, un modèle de langage à usage général récemment publié, a considérablement amélioré les performances dans les tests de mesures de risque biologique (également présentés dans la figure 2.3) et le raisonnement scientifique général par rapport aux modèles de pointe précédents (2*). Les efforts visant à formuler les meilleures pratiques en matière de biosécurité ont progressé, le Frontier Model Forum et le AI x Bio Global Forum facilitant les discussions sur l'évaluation et l'atténuation des risques pour ces modèles (388, 389).

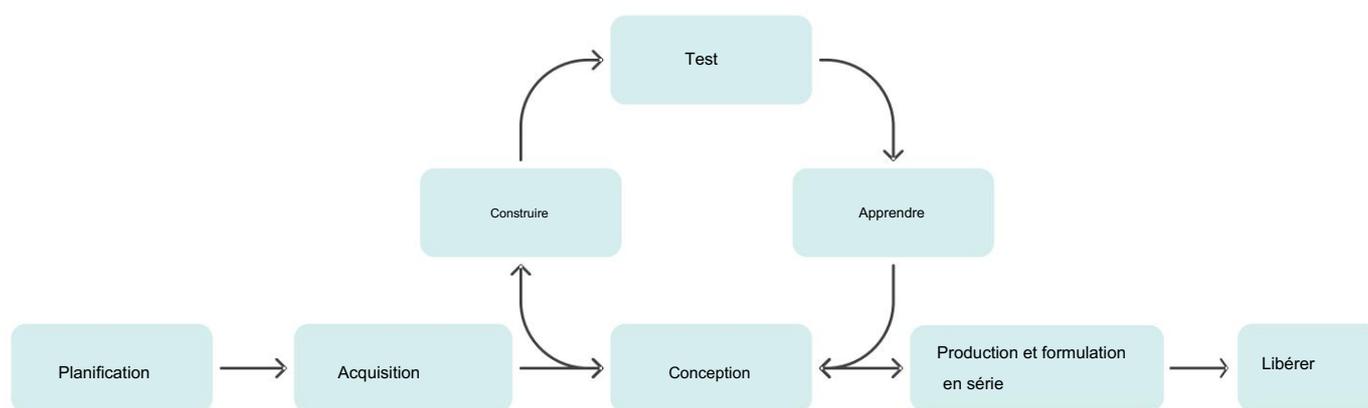


Figure 2.4 : Aperçu d'un processus de développement de produits chimiques et biologiques typique, qui est parallèle au processus utilisé pour créer des armes chimiques et biologiques. Les LLM peuvent aider aux étapes de planification et d'acquisition, conseiller sur la réalisation de travaux de laboratoire pour construire et tester une conception, et aider à planifier la mise sur le marché ou la livraison efficace d'un produit. Les agents d'IA, les plateformes robotiques et les outils de conception biologique ou chimique (à usage général ou spécialisés) peuvent aider à la conception, à la construction, aux tests et au perfectionnement des agents pathogènes et des toxines. L'IA spécialisée peut aider à la production et à la formulation de masse.

Source : Rapport international sur la sécurité de l'IA.

Les LLM peuvent désormais fournir des plans détaillés, étape par étape, pour la création d'armes chimiques et biologiques, améliorant ainsi les plans rédigés par des personnes titulaires d'un doctorat pertinent. Bien que les informations sur la manière de créer des menaces chimiques et biologiques soient accessibles depuis longtemps en raison de leur nature à double usage, les tests des LLM montrent qu'ils aident les novices à synthétiser ces informations, leur permettant d'élaborer des plans plus rapidement qu'avec Internet seul (391) (pour les phases « Planification » et « Diffusion » dans la figure 2.4). Ces capacités réduisent les obstacles à l'accès des personnes à des informations scientifiques complexes, ce qui peut probablement apporter de vastes avantages, mais peut également réduire les obstacles à l'utilisation abusive de ces informations. GPT-4, publié en 2023, a répondu correctement à 60 à 75 % des questions relatives aux armes biologiques (392), mais une gamme de modèles testés n'a apporté aucune amélioration significative par rapport aux plans d'armes biologiques développés en utilisant uniquement Internet (37*, 393, 394*). Cependant, le récent modèle o1 produit des plans jugés supérieurs aux plans générés par des experts titulaires d'un doctorat dans 72 % des cas et fournit des détails que les évaluateurs experts n'ont pas pu trouver en ligne (2*). OpenAI a conclu que ses modèles o1 pourraient aider de manière significative les experts dans la planification opérationnelle de la reproduction de menaces biologiques connues, ce qui a conduit OpenAI à augmenter son évaluation des risques biologiques de « faible » à « moyen ». Cependant, OpenAI n'a pas évalué l'utilité des modèles pour les novices (2*), soulignant la nécessité de davantage de recherches. Le développement et le déploiement réussis d'armes biologiques nécessitent toujours une expertise, des matériaux et un travail physique qualifiés importants (380, 381), ce qui signifie que même si un novice a un plan bien formulé, cela ne signifie pas qu'il pourrait le mener à bien.

Les preuves montrent que l'IA à usage général peut, dans certains cas, indiquer aux utilisateurs comment acquérir des agents biologiques et chimiques dangereux en contournant les contrôles traditionnels.

L'accès restreint aux matières dangereuses et à certains de leurs précurseurs constitue une défense essentielle contre les menaces biologiques et chimiques (phase « Acquisition » dans la figure 2.4).

Cependant, il arrive parfois que des agents biologiques soient issus de la nature ou synthétisés à partir de l'ADN, et des chimistes qualifiés peuvent identifier des voies alternatives pour créer certaines armes chimiques, en contournant les contrôles. L'IA à usage général peut aider à identifier ces alternatives

voies d'acquisition, réduisant les barrières à l'accès et créant des risques d'accidents ou de mauvaise utilisation (120). Plusieurs études suggèrent que l'IA pourrait également saper les contrôles existants sur l'accès aux séquences d'ADN à risque. De nombreux fournisseurs commerciaux d'ADN examinent leurs commandes pour détecter toute similitude avec des dangers biologiques connus afin de se conformer aux contrôles réglementaires et d'empêcher l'utilisation abusive de ces matériaux. Cependant, les LLM peuvent guider les clients vers l'achat d'ADN auprès de fournisseurs qui ne procèdent pas à des contrôles, ou suggérer des méthodes pour tromper les logiciels de dépistage (391). En outre, une étude récente a révélé que certains logiciels de dépistage ne parviennent pas à détecter une grande partie de l'ADN qui est conçu par des outils d'IA spécialisés pour fonctionner de la même manière que ces dangers, mais qui semble différent. Heureusement, la même étude a révélé qu'il est possible de mettre à jour les systèmes actuels pour détecter environ 97 % de ces conceptions (395).

Français La capacité de l'IA à concevoir des traitements médicaux hautement ciblés a considérablement augmenté depuis le rapport intermédiaire, et les interfaces de chat élargissent l'accès, augmentant également le risque de création de toxines plus puissantes (384). Les outils d'IA spécialisés et polyvalents peuvent désormais concevoir des molécules thérapeutiques candidates pour des maladies complexes telles que le cancer, les maladies auto-immunes et les affections neurologiques (la phase de « Conception » dans la figure 2.4) (396). Par exemple, AlphaProteo peut concevoir des protéines qui se fixent aux cibles jusqu'à 300 fois plus fortement que les alternatives existantes, ce qui les rend potentiellement plus efficaces à des doses plus faibles (387*). Cependant, le ciblage précis de ces systèmes pourrait également être utilisé à des fins malveillantes (397). Les outils de conception chimique pilotés par l'IA destinés à réduire la toxicité ont été réutilisés dans des études de recherche pour l'augmenter, aidant potentiellement à la conception d'armes chimiques (398), et certains outils ont été spécifiquement conçus pour la création de toxines (399). L'accès aux outils de conception spécialisés varie : certains sont réservés à des partenaires de confiance (387*), tandis que d'autres sont ouverts et peuvent donc être utilisés par n'importe qui (399). Bien que nombre de ces outils soient trop complexes pour être utilisés par des novices, des chatbots et des agents d'IA sont intégrés à certains outils de conception (400*, 401*, 402), permettant aux utilisateurs de demander des conceptions en langage clair. Aujourd'hui, cette intégration nécessite encore des connaissances techniques pour une utilisation efficace. Une évaluation directe des risques que ces outils représentent pour le développement d'armes toxiques risque d'être limitée dans les pays qui adhèrent aux engagements des traités internationaux (403).

Français L'IA à usage général améliore la capacité des chercheurs à prédire les propriétés importantes des agents pathogènes, ce qui peut aider à la fois à la conception d'armes biologiques et de contre-mesures. Des outils d'IA sont en cours de développement pour prédire les nouveaux variants de virus avant leur apparition et pour évaluer des propriétés telles que leur capacité à infecter les humains (404, 405) et à échapper à la détection immunitaire (390). Ces avancées permettent potentiellement le développement proactif de vaccins pour les variants de virus à haut risque qui n'ont pas encore émergé, ou la conception malveillante de virus capables de contourner l'immunité existante dans la population (390, 406) (la phase de « Conception » dans la Figure 2.4). Les modèles d'IA à usage général formés sur des données biologiques commencent à étayer ces applications spécialisées. Par exemple, l'outil EVEscape, qui s'appuie sur un modèle de base de l'ADN (407) et s'appuie sur des prédictions de structure protéique qui sont de plus en plus générées à l'aide de l'IA, a prédit 66 % des variants du SARS-CoV-2 (coronavirus) qui sont ensuite devenus dominants, dépassant de loin les modèles précédents (17 % de réussite) (voir Figure 2.3) (390). Les outils permettant de concevoir des virus qui échappent au système immunitaire humain et ciblent des cellules spécifiques sont utiles pour les applications de thérapie génique, mais ils présentent également des risques de double usage (408), comme l'amélioration des armes biologiques (382,

384) ou en ciblant des populations spécifiques (384). De simples modifications apportées aux agents pathogènes existants pourraient considérablement accroître leur risque. Par exemple, dans le cadre de recherches menées sans recours à l'IA, des chercheurs ont modifié des virus de la grippe aviaire qui peuvent être mortels pour les humains pour se propager par des gouttelettes en suspension dans l'air (409). En effet, certains experts pensent que les maladies créées par génie génétique pourraient être bien pires que celles qui surviennent naturellement (384, 410). L'étude de la capacité des systèmes d'IA à créer des agents pathogènes plus dangereux peut être limitée à la fois par les engagements des traités internationaux et par le risque de libération accidentelle de variantes d'agents pathogènes testées.

L'IA à usage général peut aider à planifier et à guider le travail en laboratoire, mais elle omet souvent des informations de sécurité essentielles, et les tests de réussite dans le monde réel utilisant cette assistance n'ont pas été publiés. Au moment du rapport intermédiaire (mai 2024), les plans de laboratoire générés par le LLM n'ont montré aucune amélioration significative par rapport à ceux compilés à partir d'Internet (393) (la phase « Build » dans la figure 2.4). Cependant, le modèle o1 a produit des instructions de laboratoire qui ont été préférées à celles rédigées par un doctorat dans 80 % des cas (contre 55 % pour le GPT-4), avec sa précision dans l'identification des erreurs dans les plans de laboratoire augmentant de 57 % à 73 % (2*). Malgré cela, l'omission de détails de sécurité cruciaux - comme le moment où une expérience produirait un intermédiaire explosif ou le moment où un équipement de protection doit être porté pour accomplir une tâche - reste une préoccupation importante qui pourrait entraîner des accidents graves (2*).

Les évaluations sur la manière dont les novices effectuent les travaux de laboratoire dans le cadre d'un enseignement LLM n'ont pas encore été rendues publiques, ce qui fait de la crédibilité de ces risques un sujet de débat substantiel.

L'automatisation des laboratoires et de la conception accélère le perfectionnement des conceptions biologiques. Cela réduit potentiellement les obstacles aux produits conçus par l'IA (y compris les armes), mais une mise en œuvre limitée complique l'évaluation des risques. Les conceptions biochimiques passent souvent par des cycles « Conception-Construction-Test-Apprentissage » (DBTL) pour tester et améliorer les conceptions initiales prometteuses (illustrées dans la figure 2.4). Les outils pilotés par l'IA automatisent ces cycles pour de meilleurs résultats en moins de temps (411, 412, 413, 414). Les « laboratoires autonomes » ou « scientifiques robots », un domaine naissant du développement scientifique, peuvent effectuer ces cycles sans intervention humaine (412, 415, 416) : un exemple a réalisé 20 cycles d'amélioration de conception en deux mois – soit 1 à 2 semaines de production ininterrompue – contre 6 à 12 mois manuellement (417). Les agents d'IA devraient jouer un rôle croissant dans ce processus (402, 415), et des études suggèrent que les systèmes robotiques peuvent capturer numériquement certaines des compétences motrices fines requises pour exécuter avec succès des expériences que les novices acquièrent traditionnellement au fil des années d'expérience pratique (384, 418). Si des compétences expérimentales complexes sont capturées par des plateformes robotiques, alors des capacités biologiques avancées deviendraient plus accessibles aux acteurs ayant moins de compétences techniques, mais cela n'a pas été systématiquement testé pour les compétences de laboratoire requises pour le développement d'armes biologiques. L'automatisation complète du travail de laboratoire reste un défi, par exemple en raison de pannes de machines (416).

Les applications de l'IA en biotechnologie abaissent certains obstacles à la militarisation et à la livraison d'agents chimiques et biologiques, mais ces étapes restent techniquement complexes. Des défis tels que la production de masse, la stabilisation et la dispersion efficace ont entraîné des échecs dans les programmes d'armement parrainés par les États (380, 381) et les thérapies de stade avancé (419, 420) (les phases « Production de masse et formulation » et « Libération » de la figure 2.4). Les modèles de base formés sur des données protéiques ont amélioré l'efficacité des fonctions protéiques (jusqu'à 60 %), ce qui réduit la quantité de produit nécessaire et améliore le rendement.

Français de 4x et a amélioré la stabilité des matériaux de 20 % (421), permettant une meilleure production de produits qui pourraient être utilisés comme armes ou produits thérapeutiques. Cependant, la production de masse d'organismes entiers reste difficile et les applications d'IA qui tentent de faciliter ce processus restent limitées dans leurs capacités (422). Des modèles d'IA simples peuvent également offrir un support de base pour la formulation de méthodes d'administration, telles que les poudres et les aérosols (423), un processus qui a été considéré comme un obstacle majeur au succès (tant dans le développement d'armes que de produits thérapeutiques) (424). Bien qu'un LLM récemment développé ait obtenu respectivement 100 % et 80 % de succès dans les tâches simulées de production de masse et de livraison (2*), les détails de ces tests ne sont pas disponibles, il n'est donc pas clair dans quelle mesure ils capturent les défis pratiques associés à ces étapes.

Français Les principales lacunes en matière de données probantes comprennent le manque de transparence et de cohérence dans les évaluations de sécurité et les difficultés à mesurer les capacités de conception biologique. Les évaluations de l'IA peuvent définir les « novices » comme des membres du public ou des personnes titulaires d'une licence dans un domaine spécifique, tandis que les « experts » peuvent être des personnes titulaires d'un doctorat dans une discipline pertinente ou des personnes ayant des décennies d'expérience dans un domaine spécialisé. Dans certaines évaluations, ces termes ne sont pas définis, ce qui rend difficile d'évaluer dans quelle mesure les modèles améliorent les capacités humaines et combien de personnes seraient capables d'exploiter efficacement ces capacités. Alors que de nombreuses études ont exploré le rôle de l'IA dans le développement des armes biologiques (2*, 51*, 318*, 393, 425, 426*), l'évaluation de la biosécurité de l'IA est encore un domaine naissant, avec peu de références standardisées ou d'évaluations des risques, ce qui rend difficile la comparaison des capacités et la mesure du risque créé par un nouvel outil par rapport aux technologies préexistantes (appelé « risque marginal »). L'évaluation des capacités de conception de protéines et de produits chimiques de l'IA est particulièrement difficile, car elle nécessite un processus coûteux de conception et de test des modèles. Il est peu probable que les informations clés sur les risques soient rendues publiques en raison des accords de confidentialité et des inquiétudes concernant le potentiel de sensibilisation à des pistes plus prometteuses pour militariser la biologie (427). Un dernier défi consiste à évaluer le risque global lié au développement et au déploiement d'armes biologiques et chimiques, plutôt qu'à évaluer les outils et les capacités de manière isolée.

Des efforts visant à limiter le potentiel d'utilisation abusive des systèmes d'IA à usage général formés sur des données biologiques et chimiques sont en cours, mais restent rares et sous-développés par rapport à ceux des LLM.

Français Les mesures de protection conçues pour d'autres modèles d'IA ne s'appliquent pas directement à ceux formés sur des données biologiques ou chimiques (383). Les défis liés au contrôle des résultats à risque sont doubles : 1) il existe une large gamme de résultats potentiellement dangereux provenant des outils de conception biochimique, qui ne peuvent pas être facilement définis par un filtre de contenu, et 2) les résultats bénéfiques de l'IA pour les produits thérapeutiques se chevauchent fortement (ou complètement) avec ces résultats à risque, mêlant étroitement les risques aux avantages. Alors que la communauté de conception des protéines a publié une déclaration générale sur l'utilisation responsable, des plans de mise en œuvre concrets font actuellement défaut (428). Des techniques d'atténuation des risques pour ces modèles ont été proposées mais ont jusqu'à présent reçu un développement et des tests limités (429). Cependant, certains développeurs d'IA ont exclu les données sur les agents pathogènes (386*, 430) ou restreint l'accès aux outils à haut risque (387*, 431) pour réduire les risques. Les efforts visant à empêcher les modèles d'IA à usage général de fournir des résultats à double usage sont encore compliqués par les fortes pressions de la communauté pour publier des modèles sous des licences open source et open-source (432), ce qui signifie qu'ils peuvent être téléchargés et adaptés par n'importe qui et à n'importe quelle fin (voir 2.4. Im

Par exemple, les modèles d'IA à [usage général](#) initialement formés sans séquences virales dangereuses ont ensuite été affinés avec ces données pour des applications bénéfiques (433, 434).

Français Les décideurs politiques sont confrontés à des défis majeurs pour équilibrer les avantages et les risques des capacités, en particulier lorsqu'ils fixent des limites pour une surveillance renforcée. Les modèles d'IA à usage général formés sur des données biologiques et chimiques sont souvent de pondération ouverte et moins gourmands en calcul, ce qui rend les garanties difficiles à appliquer (384), comme [indiqué dans 2.4. Impact des modèles d'IA à usage général de pondération ouverte sur les risques liés à l'IA](#) (435). Les pays qui ont signé la Convention sur les armes chimiques (CAC) et la Convention sur les armes biologiques et à toxines (BTWC) sont tenus d'empêcher le développement et l'utilisation de toute arme chimique et biologique, mais les évaluations des risques liés à l'IA se concentrent souvent uniquement sur les risques à conséquences élevées, tels que les pandémies, et négligent les risques de prolifération des armes chimiques et des toxines (318*, 410). Les décideurs politiques sont confrontés à un défi pour déterminer quelles capacités justifient des réglementations plus strictes tout en soutenant la recherche bénéfique, ce qui inclut le développement de protections contre les risques décrits dans cette section. L'évaluation de ces risques est encore compliquée par le fait que les preuves clés sont souvent des informations classifiées (427).

Les progrès dans la conception biologique se succèdent rapidement, ce qui crée une incertitude marquée quant aux capacités et aux risques futurs. Il sera essentiel de surveiller l'adoption, le succès et la sophistication de l'IA à chaque étape du processus de développement des produits biotechnologiques pour comprendre leur impact sur les programmes de biotechnologie et d'armes biologiques et la capacité des décideurs politiques à développer des mesures de prévention et de protection contre ces risques. Si une capacité transformatrice et dangereuse associée à un outil d'IA déjà commercialisé est annoncée, il se peut que peu de mesures puissent être prises pour faire face au risque.

L'élaboration d'une méthodologie d'évaluation des risques plus complète permettrait de déclencher des mesures d'atténuation avant que des risques graves ne se matérialisent, de réduire le risque que des mesures d'atténuation inutiles soient prises et permettent ainsi de bénéficier des avantages substantiels de la technologie de l'IA à usage général.

Pour les pratiques de gestion des risques liées à la science à double usage, voir :

- [3.3. Identification et évaluation des risques](#)
- [3.4.1. Former des modèles plus fiables](#)
- [3.4.2. Suivi et intervention](#)

2.2. Risques liés aux dysfonctionnements

2.2.1. Problèmes de fiabilité

INFORMATIONS CLÉS

- Le recours à des produits d'IA à usage général qui ne remplissent pas leur fonction prévue peut entraîner des dommages. Par exemple, les systèmes d'IA à usage général peuvent inventer des faits (« hallucinations »), générer un code informatique erroné ou fournir des informations médicales inexactes. Cela peut entraîner des dommages physiques et psychologiques pour les consommateurs et des dommages réputationnels, financiers et juridiques pour les individus et les organisations. • De tels problèmes de fiabilité surviennent en raison de lacunes techniques ou d'idées fausses sur les capacités et les limites de la technologie. Par exemple, les problèmes de fiabilité peuvent provenir de défis techniques tels que des hallucinations, ou du fait que les utilisateurs appliquent les systèmes à des tâches inadaptées. Les garde-fous existants pour contenir et atténuer les problèmes de fiabilité ne sont pas infaillibles.
- En raison des nombreuses utilisations potentielles de l'IA à usage général, les problèmes de fiabilité sont difficiles à prévoir. Les évaluations préalables à la sortie ne prennent pas en compte les problèmes de fiabilité qui ne se manifestent que dans le cadre d'une utilisation réelle. De plus, les techniques existantes pour mesurer les problèmes de fiabilité ne sont pas robustes, ce qui signifie qu'il n'est pas encore possible d'évaluer de manière fiable les techniques de prévention et d'atténuation.
- Les chercheurs tentent de développer des techniques de mesure et d'atténuation plus utiles, notamment pour remédier aux lacunes techniques. Depuis la publication du rapport intermédiaire (mai 2024), les mesures et les stratégies d'atténuation visant à résoudre les problèmes de fiabilité de l'IA à usage général se sont développées.
- L'un des principaux défis pour les décideurs politiques est le manque de pratiques normalisées pour prédire, identifier et atténuer les problèmes de fiabilité. Une gestion des risques sous-développée rend difficile la vérification des déclarations des développeurs sur les fonctionnalités d'IA à usage général. Les décideurs politiques sont également confrontés au défi de trouver un équilibre entre la promotion de l'innovation et la prévention d'une dépendance excessive à l'égard de l'IA.

Définitions clés

- **Fiabilité** : capacité d'un système d'IA à exécuter systématiquement sa fonction prévue.
- **Confabulations ou hallucinations** : informations inexactes ou trompeuses générées par un système d'IA, par exemple des faits ou des citations erronés.

L'IA à usage général peut être confrontée à des problèmes de fiabilité – parfois avec des conséquences dangereuses – qui ont un impact sur les personnes, les organisations et les systèmes sociaux. Les principales catégories de problèmes de fiabilité de l'IA à usage général comprennent (voir le tableau 2.2) :

- Confabulations ou hallucinations (101), c'est-à-dire contenus inexacts ou trompeurs.
- Échecs dans la réalisation de raisonnements et d'inférences de bon sens (436).
- Manque de reflet des connaissances et de la compréhension contextuellement pertinentes, à jour et impartiales (437, 438).

Les cas de défaillance de fiabilité peuvent créer des risques (439), tels que des dommages physiques ou psychologiques aux individus, des dommages à la réputation, juridiques et financiers aux organisations, et des informations erronées ayant un impact sur les processus de gouvernance.

Les exemples de problèmes de fiabilité de l'IA à usage général vont de la génération de code informatique erroné à la citation de précédents inexistant dans des mémoires juridiques. Par exemple, en génie logiciel, les LLM peuvent automatiser la génération de code informatique et aider les utilisateurs à réécrire, tester ou déboguer du code informatique (440*, 441). Cependant, les LLM ne fonctionnent souvent pas comme prévu (122, 442, 443).

Le code généré par LLM peut introduire des bugs (443), ainsi que des modifications déroutantes ou trompeuses (442).

Ces outils pourraient s'avérer utiles pour guider les programmeurs novices dans l'automatisation de certaines parties de leur flux de travail (441). Une étude de 2022 a révélé que le code des programmeurs qui utilisaient l'IA présentait davantage de vulnérabilités en matière de sécurité, et les utilisateurs n'en étaient pas conscients (444), bien que les modèles se soient considérablement améliorés depuis lors. À titre d'exemple, le modèle GPT-4 a réussi « un examen du barreau [juridique] simulé avec un score autour des 10 % des meilleurs candidats » (147*). La confiance dans ce résultat a conduit certains avocats à adopter la technologie dans leurs flux de travail professionnels (445). Dans des circonstances différentes, cependant, par exemple lorsque les paramètres de passation de l'examen étaient différents, ou lorsqu'on le comparait aux candidats au barreau qui avaient réussi l'examen la première fois qu'ils l'avaient passé (par opposition aux candidats qui avaient répété l'examen), le modèle a obtenu des performances nettement inférieures (446). Les avocats qui ont utilisé le modèle dans leur pratique juridique sans surveillance adéquate ont été confrontés à des conséquences professionnelles pour les erreurs produites par ces modèles (447). Des malentendus similaires concernant la fiabilité du modèle s'appliquent dans le contexte médical (448) : les modèles ont passé avec succès des tests médicaux (147*, 449) et sont censés posséder des connaissances cliniques fiables, mais l'utilisation dans le monde réel et les réévaluations nuancées révèlent des limites (450).

Les principales causes des problèmes de fiabilité de l'IA à usage général sont 1. les limitations technologiques et 2. les idées fausses sur les capacités des modèles (456). Certaines des principales limitations technologiques de l'IA à usage général sont répertoriées dans le tableau 2.2. Les idées fausses sur la technologie et l'absence de garde-fous de sécurité adéquats peuvent conduire à une dépendance excessive et à ce que les utilisateurs appliquent les systèmes à des tâches impossibles et pratiquement difficiles que l'IA à usage général n'est pas capable d'accomplir (456). Les limitations et les idées fausses sont exacerbées par l'incitation à publier des modèles et des produits d'IA à usage général avant qu'ils n'aient été correctement évalués et que leurs capacités et limitations n'aient fait l'objet de recherches scientifiques.

Type de problèmes de fiabilité	Exemples
Confabulations ou hallucinations	<ul style="list-style-type: none"> • Citer des précédents inexistants dans des mémoires juridiques (451) • Citant des politiques de tarifs réduits inexistantes pour les passagers en deuil (452)
Les échecs du raisonnement de bon sens	<ul style="list-style-type: none"> • Ne pas effectuer de calculs mathématiques de base (453*) • Ne pas déduire de relations causales de base (454)
Les échecs de la connaissance contextuelle	<ul style="list-style-type: none"> • Fournir des informations médicales inexactes (448) • Fournir des informations obsolètes sur des événements (455)

Tableau 2.2 : L'IA à usage général peut présenter divers problèmes de fiabilité.

Étant donné la nature polyvalente et l'utilisation généralisée de l'IA à usage général, tous les problèmes de fiabilité ne peuvent pas être prévus et suivis. Plusieurs mécanismes existent pour prévoir et suivre les problèmes de fiabilité dans l'IA à usage général. Il s'agit notamment d'évaluations visant à évaluer la prévalence de divers problèmes de fiabilité avant la sortie du produit (457, 458) et de la tenue à jour de référentiels d'incidents d'IA (tels que l'AI Incidents Monitor (AIM) de l'OCDE (459)) après la sortie afin d'éviter des incidents similaires à l'avenir.

Toutefois, étant donné la nature polyvalente de la technologie et ses cas d'utilisation toujours croissants dans de nouveaux domaines, il n'est pas garanti que de tels mécanismes permettent de détecter tous les risques possibles.

Les garde-fous existants pour contenir et atténuer les problèmes de fiabilité ne sont pas infaillibles (460). Par exemple, bien que des travaux récents aient proposé des méthodes pour atténuer les hallucinations (461), il n'existe aucune preuve solide de l'efficacité de ces méthodes, et il n'existe pas de méthodes infaillibles pour atténuer les hallucinations. Pour promouvoir la fiabilité de l'IA à usage général, les évaluateurs doivent évaluer les systèmes de manière rigoureuse avant leur publication, communiquer de manière précise et accessible sur les résultats et sur la manière dont ils doivent ou ne doivent pas être interprétés par les utilisateurs, et spécifier les utilisations prévues des systèmes (et les utilisations qui ne sont pas prévues).

Depuis la publication du rapport intermédiaire, le référentiel de mesures et de stratégies d'atténuation des problèmes de fiabilité de l'IA à usage général n'a cessé de s'élargir. Par exemple, un consortium de chercheurs, d'ingénieurs et de praticiens de l'industrie et du monde universitaire a développé un « référentiel de sécurité de l'IA » (457), qui vise à évaluer la sécurité spécifique à chaque cas d'utilisation.

Les risques des systèmes d'IA basés sur LLM en proposant une approche raisonnée pour la construction de tests de référence et une plateforme ouverte pour tester un large éventail de dangers. COMPL-AI est un autre cadre d'évaluation open source récemment publié pour les modèles d'IA génératifs (462).

L'objectif de cette étude est d'évaluer la conformité des modèles d'IA aux exigences de la loi européenne sur l'IA en termes de robustesse, de confidentialité, de droits d'auteur et au-delà (458). Les chercheurs ont continué à proposer de nouveaux critères de référence (par exemple pour le raisonnement causal (454) ou le raisonnement juridique (463)) et ont étudié les lacunes des critères de référence existants (178, 464).

La principale lacune en matière de données probantes sur les problèmes de fiabilité de l'IA à usage général concerne l'efficacité des mécanismes existants pour atténuer ces problèmes. Par exemple, la conception d'évaluations fiables et reproductibles des capacités, des limites et des échecs de l'IA à usage général avant,

pendant et après le déploiement reste un défi majeur (465). De plus, certains problèmes de fiabilité (par exemple, le recours à des informations obsolètes (455)) peuvent ne se manifester que dans le cadre d'une utilisation réelle, ce qui rend les évaluations préalables à la publication inadéquates. Le développement et la maintenance de bancs d'essai collaboratifs en évolution dynamique pour évaluer les fonctionnalités de l'IA à usage général peuvent être une piste pour combler ces lacunes. Parmi les autres lacunes critiques, citons le manque de bonnes pratiques pour une publication responsable des produits.

Les décideurs politiques qui souhaitent promouvoir la fiabilité de l'IA à usage général sont confrontés à plusieurs compromis et défis. Compte tenu de l'utilisation généralisée de la technologie, il est important que les produits et services d'IA à usage général fonctionnent comme prévu (456). Cependant, les normes et les meilleures pratiques requises n'ont pas encore été correctement établies (457, 465, 466). En outre, il est difficile de garantir le respect des meilleures pratiques existantes en l'absence d'incitations, d'organismes d'évaluation de la conformité et d'évaluateurs experts possédant les compétences sociotechniques requises (467). L'un des principaux problèmes est l'incertitude entourant l'efficacité des mécanismes existants pour prédire et atténuer les risques d'échec. L'absence d'exigences normalisées pour évaluer et documenter les capacités et les limites des modèles rend difficile la vérification des déclarations des développeurs sur la fiabilité de l'IA à usage général – une condition préalable à l'élaboration efficace des politiques en matière d'IA (468). Un autre défi consiste à équilibrer la nécessité de promouvoir l'innovation et la compétitivité économique tout en décourageant les déclarations non fondées et la dépendance excessive à l'égard de la technologie. La lutte contre la dépendance excessive nécessite d'évaluer et d'améliorer l'état actuel des connaissances en IA parmi les utilisateurs et les consommateurs de la technologie. Les outils et les idées provenant d'industries critiques en matière de sécurité plus matures peuvent offrir des conseils utiles pour relever les défis mentionnés ci-dessus, mais le rythme des progrès technologiques peut compliquer ces efforts.

Pour les pratiques de gestion des risques liées aux problèmes de fiabilité, voir :

- [3.3. Identification et évaluation des risques](#)
- [3.4.1. Former des modèles plus fiables](#)
- [3.4.2. Suivi et intervention](#)
- [3.4.3. Moyens techniques de protection de la vie privée](#)

2.2.2. Biais

INFORMATIONS CLÉS

- Les systèmes d'IA à usage général peuvent amplifier les préjugés sociaux et politiques, provoquant des préjudice. Ils affichent souvent des préjugés liés à la race, au sexe, à la culture, à l'âge, au handicap, à l'opinion politique ou à d'autres aspects de l'identité humaine. Cela peut conduire à des résultats discriminatoires, notamment une répartition inégale des ressources, le renforcement des stéréotypes et la négligence systématique de certains groupes ou points de vue.
- Les biais dans l'IA ont de nombreuses sources, comme des données de formation médiocres et des choix de conception du système. L'IA à usage général est principalement formée sur des ensembles de données linguistiques et d'images qui Les cultures anglophones et occidentales sont représentées de manière disproportionnée. Cela contribue à un résultat biaisé. Certains choix de conception, tels que les techniques de filtrage de contenu utilisées pour aligner les systèmes sur des visions du monde particulières, peuvent également contribuer à un résultat biaisé.
- Les mesures d'atténuation techniques ont conduit à des améliorations substantielles, mais ne fonctionnent pas toujours. Les chercheurs ont fait des progrès significatifs dans la lutte contre les biais dans l'IA à usage général, mais plusieurs problèmes restent encore non résolus. Par exemple, la frontière entre les stéréotypes nuisibles et la connaissance utile et précise du monde peut être difficile à tracer, et la perception des biais peut varier en fonction des contextes culturels, des contextes sociaux et des cas d'utilisation.
- Depuis la publication du rapport intermédiaire (mai 2024), la recherche a révélé de nouvelles des types plus subtils de biais de l'IA. Par exemple, des travaux récents ont montré que l'IA à usage général peut générer des résultats biaisés selon que l'utilisateur interagit avec l'IA dans un certain dialecte.
- Les décideurs politiques sont confrontés à des compromis liés aux biais de l'IA. Il existe de nombreux domaines, tels que le droit La prise de décision, pour laquelle l'IA polyvalente peut en principe s'avérer très utile, n'est toutefois pas toujours fiable, ce qui peut entraîner des risques de discrimination. Les décideurs politiques doivent peser les compromis fondamentaux entre des priorités concurrentes telles que l'équité, l'exactitude et la confidentialité, en particulier lorsqu'ils réglementent des applications à enjeux élevés.

Définitions clés

- **Biais** : erreurs systématiques dans les systèmes algorithmiques qui favorisent certains groupes ou visions du monde et créent souvent des résultats injustes pour certaines personnes. Les biais peuvent avoir de multiples sources, notamment des erreurs de conception algorithmique, des ensembles de données non représentatifs ou autrement erronés, ou des inégalités sociales préexistantes.
- **Discrimination** : Le traitement injuste d'individus ou de groupes en fonction de leurs attributs, tels comme la race, le sexe, l'âge, la religion ou d'autres caractéristiques protégées.
- **Collecte et prétraitement des données** : une étape du développement de l'IA au cours de laquelle les développeurs et Les travailleurs des données collectent, nettoient, étiquettent, normalisent et transforment les données de formation brutes dans un format à partir duquel le modèle peut apprendre efficacement.

- Apprentissage par renforcement à partir du retour d'information humain (RLHF) : une technique d'apprentissage automatique dans laquelle un modèle d'IA est affiné en utilisant des évaluations ou des préférences fournies par l'homme comme signal de récompense, permettant au système d'apprendre et d'ajuster son comportement pour mieux s'aligner sur les valeurs et les intentions humaines grâce à un entraînement itératif.
- IA explicable (XAI) : un programme de recherche visant à créer des systèmes d'IA qui fournissent des informations claires et des explications compréhensibles de leurs décisions, permettant aux utilisateurs de comprendre comment et pourquoi des résultats spécifiques sont générés.

Il existe plusieurs cas bien documentés de systèmes d'IA, polyvalents ou non, amplifiant les préjugés sociaux ou politiques. Cela peut, par exemple, prendre la forme de résultats discriminatoires fondés sur la race, le sexe, l'âge et le handicap, avec des effets néfastes dans des domaines tels que la santé, l'éducation et la finance. Dans les systèmes d'IA restreints, des préjugés raciaux ont été documentés dans les algorithmes de reconnaissance faciale (469), les prédictions de récidive (470, 471) et les outils de santé, qui sous-estiment les besoins des patients issus de milieux raciaux et ethniques marginalisés (472). L'IA polyvalente présente également de tels préjugés, par exemple des préjugés raciaux dans les contextes cliniques (448, 473), et il a été démontré que les générateurs d'images reproduisent des stéréotypes dans les professions (474, 475, 476). Les chercheurs ont également constaté que les modèles de génération d'images reproduisent de manière excessive les stéréotypes de genre dans des professions telles que pilote (homme) ou coiffeur (femme) et surreprésentent les personnes blanches dans tous les domaines, à l'exception des professions telles que pasteur ou rappeur (476).

Dans de nombreux cas, les biais d'IA surviennent lorsque certains groupes sont sous-représentés dans les données d'entraînement ou représentés d'une manière qui imite les stéréotypes sociétaux. Il a été démontré que les ensembles de données utilisés pour entraîner les modèles d'IA sous-représentent divers groupes de personnes, par exemple les personnes d'un certain âge, d'une certaine race, d'un certain sexe et d'un certain statut de handicap (477, 478) et sont limités en termes de diversité géographique (479*, 480). Les ensembles de données d'entraînement sont également très susceptibles d'être en anglais et de représenter les cultures occidentales (481). Ces ensembles de données sont également principalement agrégés à partir de livres numérisés et de textes en ligne, qui ne reflètent pas les traditions orales et les cultures non numérisées, potentiellement au détriment de groupes marginalisés tels que les communautés autochtones. Un tel biais de représentation peut conduire à des échecs dans la manière dont les modèles formés sur ces données sont capables de généraliser aux populations cibles (482). Par exemple, un modèle d'IA à usage général destiné à aider les femmes enceintes dans les zones rurales du Malawi ne fonctionnera pas comme prévu s'il est formé sur des données provenant de mères vivant dans les zones urbaines du Canada. En outre, les biais historiques intégrés dans les données peuvent perpétuer des injustices systémiques, telles que le financement hypothécaire injuste pour les populations minoritaires aux États-Unis (483*), conduisant potentiellement les systèmes d'IA à refléter les cultures, les langues et les visions du monde dominantes, au détriment des groupes sous-représentés dans ces systèmes (484, 485, 486, 487).

Les biais de données proviennent de facteurs historiques ainsi que de la manière dont les ensembles de données sont collectés, annotés et préparés pour l'apprentissage du modèle. Le biais de représentation se produit en raison de facteurs tels qu'une collecte et un prétraitement des données erronés, ainsi que de biais historiques tels que le racisme et le sexisme (488). En ce qui concerne la collecte de données, un biais peut émerger du choix de la source de collecte de données par le chercheur (API externes, sources de données publiques, scraping Web, etc.) (489). Au cours du processus d'étiquetage des données, un biais de mesure peut se produire lors de la sélection des étiquettes et des fonctionnalités des ensembles de données à utiliser

pour la tâche de prédiction respective, étant donné que certaines constructions abstraites comme le potentiel académique sont évaluées à l'aide de notes et de résultats aux tests (482). Dans d'autres cas, ce biais peut être exacerbé lorsque les chercheurs relèguent les tâches d'étiquetage à des annotateurs qui peuvent ne pas disposer d'un contexte culturellement pertinent pour comprendre les memes, les textes sarcastiques ou les blagues.

Les biais sont présents à différentes étapes du cycle de vie de l'apprentissage automatique, de la collecte des données au déploiement (voir le tableau 2.3). Les études sur l'IA à usage général ont de plus en plus mis en évidence les biais dans les résultats des chatbots et des générateurs d'images. À mesure que les systèmes d'IA à usage général s'intègrent progressivement dans des contextes réels, il est important de comprendre les impacts des biais de déploiement, qui peuvent se produire lorsque les systèmes d'IA sont mis en œuvre dans des contextes différents de ceux pour lesquels ils ont été conçus. Pour comprendre les limites des systèmes d'IA à usage général dans divers contextes, un certain nombre de méthodes ont été proposées pour évaluer les capacités des modèles d'IA à usage général ; cependant, celles-ci sont également sujettes à des biais. Les repères tels que la mesure de la compréhension du langage multitâche massif (MMLU), qui est un repère largement utilisé pour évaluer les capacités, sont centrés sur les États-Unis et contiennent des questions triviales et erronées (490). Bien que les travaux récents se soient concentrés sur l'atténuation des défis de ces repères (490), des recherches importantes sont nécessaires pour élargir la portée des méthodes d'évaluation afin d'inclure des contextes non occidentaux.

Les préjugés sexistes sont largement étudiés, avec des preuves détaillant leur impact sur les cas d'utilisation de l'IA à usage général et de l'IA à usage restreint. Des études empiriques ont documenté des schémas de langage sexistes et des représentations stéréotypées dans les résultats générés par l'IA à usage général (491, 492) et des résultats à prédominance masculine issus de recherches Internet neutres en termes de genre utilisant des algorithmes d'IA à usage restreint (493).

Dans le domaine de l'IA à usage général, ces problèmes se traduisent par des résultats stéréotypés, tant de la part des LLM que des générateurs d'images. Ces stéréotypes impliquent souvent des préjugés sexistes professionnels (494, 495, 496, 497).

La discrimination fondée sur l'âge dans l'IA est un domaine peu étudié par rapport à la race et au sexe, mais les premières données suggèrent que cette forme de biais de l'IA a des impacts significatifs. En 2023, les études menées lors d'une importante conférence sur l'équité, la responsabilité et la transparence (FAccT) étaient deux fois plus susceptibles d'aborder la race et le sexe que l'âge (498). De plus en plus de recherches mettent en évidence un biais lié à l'âge dans l'IA à usage général, des études antérieures l'ayant identifié dans la recherche d'emploi (499) et les prêts (500). Les LLM excluent souvent les personnes âgées dans les modèles de texte en image et génèrent des sujets de contenu biaisés liés au vieillissement (498).

Des études ont également révélé que les modèles de générateur d'images représentent en grande partie des adultes âgés de 18 à 40 ans lorsqu'aucun âge n'est spécifié, stéréotypant les adultes plus âgés dans des rôles limités (501). Une discrimination fondée sur l'âge a également été identifiée dans des LLM importants (502*, 503). Les biais dans les données de formation, où les adultes plus âgés sont sous-représentés, sont une des principales raisons de cette discrimination (504). Les résultats peuvent également être biaisés en faveur des individus plus jeunes en raison d'un biais d'incitation, l'influence involontaire des invites de saisie sur les résultats du modèle d'IA, ce qui peut conduire à des réponses biaisées ou faussées en fonction de la formulation, du contexte ou du cadrage de l'invite (501, 505).

Les préjugés liés au handicap dans l'IA sont également un domaine peu étudié, mais les recherches émergentes se concentrent sur les impacts spécifiques des systèmes d'IA à usage général sur les personnes handicapées. Les chercheurs ont montré comment les systèmes et outils d'IA à usage général peuvent discriminer les utilisateurs handicapés, par exemple

en reproduisant les stéréotypes sociétaux sur les handicaps (506) et en classant de manière inexacte les sentiments à l'égard des personnes handicapées (507). Des recherches supplémentaires ont montré les limites de ces outils pour le filtrage des CV (508) et la génération d'images (506). Les problèmes de biais liés au handicap sont également exacerbés par le manque d'ensembles de données inclusifs. Malgré des recherches croissantes sur la reconnaissance de la langue des signes, les systèmes d'IA à usage général ont des capacités de transcription limitées en raison de la rareté des ensembles de données sur la langue des signes par rapport aux langues écrites et parlées (212). La plupart des ensembles de données se concentrent sur la langue des signes américaine, ce qui limite les capacités de transcription des LLM tels que ChatGPT pour d'autres langues des signes, comme la langue des signes arabe (509). Les efforts récents pour développer des ensembles de données pour les langues des signes africaines (510) constituent un modeste pas vers une inclusion plus équitable des divers dialectes des signes.

Français Les systèmes d'IA à usage général présentent des biais politiques variés, et certaines premières preuves suggèrent que cela peut influencer les convictions politiques des utilisateurs. Des études récentes ont démontré que les systèmes d'IA à usage général peuvent être politiquement biaisés, différents systèmes favorisant différentes idéologies sur un spectre allant des opinions progressistes aux opinions centristes et conservatrices (511, 512, 513, 514, 515, 516*, 517, 518). Des études montrent également qu'un seul système d'IA à usage général peut favoriser différentes positions politiques en fonction de la langue de l'invite (519, 520) et du sujet en question (521). Par exemple, une étude a révélé qu'un système d'IA à usage général produisait des résultats plus conservateurs dans des langues souvent associées à des sociétés plus conservatrices et des résultats plus libéraux dans des langues souvent associées à des sociétés plus progressistes (520). Les biais politiques proviennent de diverses sources, notamment des données de formation qui reflètent des idéologies particulières, des modèles de réglage fin basés sur les commentaires d'évaluateurs humains biaisés et des filtres de contenu introduits par les entreprises d'IA pour exclure des résultats particuliers (520, 522). Il existe des preuves que l'interaction avec des systèmes d'IA à usage général biaisés peut affecter les opinions politiques des utilisateurs (523) et accroître la confiance dans les systèmes qui correspondent à l'idéologie de l'utilisateur (524). Cependant, des recherches supplémentaires sont nécessaires pour évaluer l'impact global de l'IA à usage général biaisée sur les opinions politiques des gens.

Les systèmes d'IA peuvent présenter des biais aggravants, où les individus ayant plusieurs identités marginalisées (par exemple une femme de couleur à faible revenu) sont confrontés à une discrimination aggravée, mais les preuves à ce sujet sont naissantes et peu concluantes. Alors que des recherches émergent sur la détection des biais aggravants dans les modèles d'IA (525, 526, 527), les progrès pour atténuer ces biais ont été plus lents (528). Des études ont montré que les modèles d'IA utilisés dans la sélection des CV et la génération de contenu d'actualité favorisent souvent les noms de femmes blanches par rapport aux noms de femmes noires (529), et que les personnes et les femmes noires sont plus sujettes à la discrimination (530). Cependant, dans certains cas, les hommes hispaniques (531) ou noirs ont obtenu les pires résultats (529). Bien que ces recherches se développent, la tendance de l'IA à usage général à afficher des biais aggravants, en particulier dans les catégories d'identité non occidentales telles que la tribu et la caste, reste globalement sous-explorée. Les modèles d'IA étant de plus en plus utilisés à l'échelle mondiale, il sera crucial de comprendre ces biais et leurs relations complexes avec la race, le sexe et d'autres identités.

Les méthodes techniques courantes de débiaisage comprennent les stratégies de prétraitement, de traitement en cours et de post-traitement (532, 533). Les « techniques de prétraitement » tentent d'éliminer le biais existant dans les données utilisées pour former les modèles d'IA. Cette classe de techniques garantit que les données sont propres et équilibrées en fonction des attributs démographiques. Les « techniques de traitement en cours » se concentrent sur la modification de l'IA

processus ou architecture de formation du modèle pour réduire les biais (voir également [3.4.1. Formation de modèles plus fiables](#) pour des méthodes similaires appliquées à une variété de problèmes). Les approches de « post-traitement » modifient les résultats de l'IA pour qu'ils soient moins biaisés (voir également [3.4.2. Suivi et intervention](#) pour des techniques similaires appliquées à une variété de problèmes). Chaque technique a ses limites ; ainsi, de nombreuses entreprises d'IA utilisent une combinaison de méthodes pour réduire progressivement les biais (30, 534*).

Une approche holistique et participative qui inclut une variété de perspectives et de parties prenantes est essentielle pour atténuer les biais. Des équipes interdisciplinaires combinant une expertise technique, juridique et sociale pour un contrôle complet des biais sont essentielles (535, 536). L'alignement des systèmes d'IA sur les valeurs sociétales est intrinsèquement difficile dans des communautés diverses, où les points de vue peuvent être contradictoires (438, 537, 538). Une représentation accrue des groupes marginalisés (539) et un dialogue participatif (538) visent à remédier aux risques liés au fait de favoriser des intérêts particuliers ; toutefois, la participation à elle seule ne suffit pas à résoudre pleinement ces conflits (540).

Il est difficile de répondre efficacement aux préoccupations en matière de discrimination, car les méthodes d'atténuation des biais ne sont pas fiables. Les problèmes d'atténuation des biais existaient avant les systèmes d'IA à usage général (541), mais les techniques actuelles de lutte contre les biais peuvent créer involontairement de nouveaux biais malgré des progrès considérables dans ce domaine. Par exemple, le RLHF introduit parfois des biais en fonction de la diversité des fournisseurs de commentaires (542). D'autres méthodes, telles que la réannotation des ensembles de données, peuvent améliorer la cohérence entre les étiquetages, mais elles sont coûteuses et prennent du temps (543, 544). Les efforts d'atténuation robustes en sont encore à leurs débuts (545). Un défi important dans l'atténuation des risques de l'IA réside également dans la définition et la mesure des résultats efficaces, en particulier pour les biais. On ne sait toujours pas comment mesurer les biais en général, comment faire la distinction entre les données qui reflètent des différences démographiques légitimes (par exemple, la prévalence des maladies par population) et les données qui perpétuent intrinsèquement les biais, et à quoi ressemble un état final idéal et mesurable. Par exemple, l'atténuation des biais à l'encontre de petits groupes ethniques ou religieux est complexe ; il est difficile de savoir quand les biais ont été suffisamment réduits, ce qui rend ces défis particulièrement prononcés dans les questions liées aux biais, bien que des écarts de mesure similaires existent pour les risques

L'évaluation de l'atténuation des biais dans les systèmes d'IA avancés repose sur des mesures quantitatives, des évaluations qualitatives et la mesure de l'impact dans le monde réel. L'objectif de ces évaluations est de mesurer le succès des techniques d'atténuation dans la réduction des biais, l'amélioration de l'équité et l'obtention de résultats équitables pour diverses populations. Ces évaluations guident également les approches d'atténuation, établissent des repères et garantissent l'alignement avec les exigences de gouvernance et/ou réglementaires (535). Les repères sont essentiels dans les domaines à enjeux élevés pour répondre aux normes juridiques et éthiques (492, 546). De plus, une surveillance continue dans le monde réel garantit que les mesures de réduction conduisent à des résultats moins biaisés dans la pratique. Par exemple, des audits réguliers des modèles d'IA utilisés dans la justice pénale peuvent vérifier que les efforts de réduction des biais restent efficaces à mesure que de nouvelles données sont introduites (547).

La réduction des biais peut entrer en conflit avec d'autres desiderata, et il peut s'avérer techniquement impossible d'atteindre une équité algorithmique complète. De nombreuses propriétés souhaitables des systèmes d'IA à usage général impliquent des compromis, tels que le compromis à quatre voies entre équité, précision, confidentialité et efficacité (548, 549*, 550, 551, 552). Les tentatives visant à garantir l'équité peuvent avoir des inconvénients. Par exemple, Gemini

généralisé des images historiquement inexactes, représentant des autochtones et des femmes de couleur comme des sénateurs américains des années 1800, ou représentant des soldats allemands de la Seconde Guerre mondiale avec des ethnies diverses. Ces images ont déformé l'histoire, peut-être en raison d'une tentative de garantir la diversité raciale dans les images générées qui n'ont pas réussi à prévoir et à s'adapter à ces cas d'utilisation spécifiques, et en raison de l'incapacité à prévoir les incitations spécifiques qui ont conduit à ce résultat. Pour faire face à ces complexités, une approche équilibrée utilisant à la fois des mesures quantitatives et qualitatives peut aider les technologues à faire des compromis éclairés. Cependant, la faisabilité technique de parvenir à une équité algorithmique complète dans les systèmes d'IA à usage général est débattue. Les résultats mathématiques indiquent qu'il peut être impossible de satisfaire simultanément tous les critères d'équité, comme le suggère un « théorème d'impossibilité » sur l'équité (550, 553, 554, 555). Même s'il existe des limites théoriques à l'équité, des solutions pratiques sont réalisables (556, 557). Certains chercheurs soutiennent que les définitions de l'équité peuvent être partiellement conciliées et que plusieurs critères d'équité peuvent être satisfaits simultanément dans une large mesure (557, 558). Des études empiriques remettent en question l'inévitabilité des compromis entre équité et précision dans les systèmes d'IA, suggérant que la réduction des biais n'entraîne souvent pas de perte significative de précision ou ne nécessite pas de méthodes complexes à mettre en œuvre (557, 558, 559).

Biais des étapes du cycle de vie	Source	Description	Exemples
Collecte de données par échantillonnage	Biais	Certaines perspectives, données démographiques ou groupes sont surreprésentés ou sous-représentés dans les données.	Un ensemble de données pour un agrégateur de nouvelles contenant principalement des sources qui favorisent une idéologie particulière, ce qui conduit à des résultats biaisés
	Sélection Biais	Seuls certains types de données ou contextes sont inclus, ce qui limite la représentativité.	Ensembles de données linguistiques qui excluent les langues non occidentales, limitant ainsi les performances du modèle dans les applications mondiales.
Données Annotation	Étiqueteuse Biais	Les antécédents, les perspectives et les préjugés culturels des annotateurs affectent leur compréhension et classification des données, influençant le processus d'étiquetage.	Les annotateurs étiquettent les discours des individus de niveau inférieur. les milieux socioéconomiques comme non professionnels ou inappropriés, conduisant à des décisions biaisées.
Conservation des données	Historique Biais	Reflétant ou perpétuant les préjugés sociétaux passés dans des œuvres sélectionnées données.	Un ensemble de données d'embauche qui favorise certaines données démographiques en fonction des pratiques d'embauche historiques, intégrant les inégalités existantes dans les modèles d'IA.
Données Prétraitement	Fonctionnalité Sélection Biais	Exclusion de fonctionnalités pertinentes d'un ensemble de données.	L'exclusion de l'âge ou du sexe comme caractéristiques des modèles de soins de santé peut avoir un impact potentiel sur la pertinence des prévisions pour ces données démographiques.

Formation de modèle	Étiquette Déséquilibre	Représentation inégale dans les données étiquetées, conduisant à des résultats de modèle biaisés.	Un modèle de classification formé sur Les images étiquetées à 80 % comme étant masculines peuvent être peu performantes lors de l'identification des images féminines.
Déploiement Contexte	Contextuel Biais	Un modèle est formé sur des données provenant de un contexte qui diffère de celui contexte d'application, conduisant à des résultats pires pour certains groupes.	Un modèle uniquement en anglais déployé dans des environnements multilingues, provoquant des interprétations erronées pour les non-anglophones utilisateurs.
Évaluation & Validation	Référence Biais	Les critères d'évaluation favorisent certains groupes ou bases de connaissances par rapport à d'autres.	Les modèles d'IA évalués principalement sur des ensembles de données centrés sur les États-Unis ne parviennent pas à se généraliser correctement dans des contextes non occidentaux.
Retour Mécanismes	Retour Biais de boucle	Les modèles apprennent à partir des commentaires biaisés des utilisateurs, renforçant ainsi les biais initiaux.	Un système de recommandation qui reçoit plus d'engagement sur certains types de contenu peut renforcer l'exposition au même contenu biaisé.

Tableau 2.3 : Les biais peuvent survenir à différentes étapes du cycle de vie de production des données pour les systèmes d'IA et avoir différentes sources, telles que des ensembles de données non représentatifs, un étiquetage biaisé ou des repères biaisés.

Depuis la publication du rapport intermédiaire, des études ont révélé de nouvelles formes plus subtiles de biais de l'IA, tandis qu'une attention accrue portée aux techniques d'atténuation et à l'explicabilité constituent des étapes importantes pour réduire les biais dans les systèmes d'IA à usage général.

- Des études récentes ont montré que les modèles linguistiques répondent différemment aux différents dialectes anglais, avec des réponses différentes à l'anglais vernaculaire afro-américain (AAVE) par rapport à l'anglais américain standard (183, 438, 491, 511, 560, 561, 562, 563, 564, 565).
Les recherches axées sur l'examen des préjugés dans les langues non occidentales ont également augmenté, démontrant l'existence de préjugés sexistes dans les modèles hindi, qui impliquent souvent des nuances subtiles (566).
Les recherches ont également exploré le « biais d'homogénéité », une forme de biais dans laquelle certains groupes sociaux sont perçus comme moins diversifiés ou plus homogènes que d'autres (567). • Les recherches sur l'atténuation des biais ont également augmenté, notamment les études sur la réduction des biais d'étiquetage (568, 569, 570). Cependant, des progrès beaucoup plus importants sont nécessaires pour comprendre dans quelle mesure ces méthodes seront efficaces pour atténuer les défis existants en matière de biais dans les systèmes du monde réel.
- Progrès en matière d'IA X : les avancées technologiques se sont de plus en plus concentrées sur l'explicabilité des LLM. Des techniques telles que les gradients intégrés et le raisonnement sur les graphes (RoG) (571, 572, 573) ont été développées pour rendre les processus de décision des modèles plus transparents. Ces méthodes pourraient faciliter la détection des biais au sein des modèles et favoriser la confiance en offrant des explications claires et interprétables des processus de prise de décision de l'IA.

Les mesures d'atténuation des biais sont souvent imparfaites, ce qui rend difficile de tirer parti des avantages de l'IA sans perpétuer divers biais. Dans certains domaines, comme la prise de décision juridique, l'IA pourrait contribuer à atténuer les biais. Cependant, les capacités actuelles de l'IA ne sont pas toujours fiables, ce qui rend difficile pour les décideurs politiques de décider si les modèles sont suffisamment sûrs pour être déployés sans menacer les droits de l'homme ou d'autres droits. En outre, l'équité n'a pas de définition universellement acceptée, sa signification variant considérablement selon les contextes culturels, sociaux et disciplinaires (574, 575, 576, 577). Les décideurs politiques devront également réfléchir aux meilleurs moyens d'impliquer les communautés les plus touchées et les plus vulnérables dans ces décisions.

À mesure que l'ampleur du déploiement de l'IA à usage général s'élargit, les difficultés à prouver les préjudices causés par la discrimination peuvent également rendre difficile l'intervention des décideurs politiques.

Pour les pratiques de gestion des risques liées aux préjugés et à la sous-représentation, voir :

- [3.3. Identification et évaluation des risques](#)
- [3.4.2. Suivi et intervention](#)

2.2.3. Perte de contrôle

INFORMATIONS CLÉS

- Les scénarios de « perte de contrôle » sont des scénarios futurs hypothétiques dans lesquels un ou plusieurs systèmes d'IA à usage général finissent par échapper au contrôle de quiconque, sans possibilité claire de reprendre le contrôle. Ces scénarios varient en gravité, mais certains experts donnent du crédit à des conséquences aussi graves que la marginalisation ou l'extinction de l'humanité.
- L'opinion des experts sur la probabilité d'une perte de contrôle varie considérablement. Certains considèrent qu'il est peu plausible, d'autres le considèrent comme probable, tandis que d'autres le considèrent comme un risque de probabilité modeste qui mérite d'être pris en compte en raison de sa gravité élevée. Les recherches empiriques et mathématiques en cours font progressivement avancer ces débats.
- Deux conditions essentielles pour les scénarios de perte de contrôle fréquemment évoqués sont a. des capacités d'IA nettement accrues et b. l'utilisation de ces capacités de manière à saper le contrôle. Premièrement, certains futurs systèmes d'IA auraient besoin de capacités spécifiques (surpassant largement celles des systèmes actuels) qui leur permettraient de saper le contrôle humain. Deuxièmement, certains systèmes d'IA devraient employer ces « capacités de sape du contrôle », soit parce qu'ils ont été intentionnellement conçus pour le faire, soit parce que des problèmes techniques produisent un comportement inattendu.
- Depuis la publication du rapport intermédiaire (mai 2024), les chercheurs ont observé des progrès modestes vers le développement de capacités de sape du contrôle. Les capacités pertinentes incluent des capacités de planification autonomes associées aux agents d'IA, des capacités de programmation plus avancées et des capacités utiles pour saper la surveillance humaine.
- La gestion d'une éventuelle perte de contrôle pourrait nécessiter une préparation importante en amont, malgré les incertitudes existantes. L'un des principaux défis pour les décideurs politiques est de se préparer à un risque dont la probabilité, la nature et le moment de sa survenance demeurent inhabituellement ambigus.

Définitions clés

- Contrôle : La capacité d'exercer une surveillance sur un système d'IA et d'ajuster ou d'arrêter son comportement si il agit de manière indésirable.
- Scénario de perte de contrôle : scénario dans lequel un ou plusieurs systèmes d'IA à usage général se mettent à fonctionner en dehors du contrôle de quiconque, sans voie claire pour reprendre le contrôle.
- Capacités de contrôle : Capacités qui, si elles étaient utilisées, permettraient à un système d'IA pour saper le contrôle humain.
- Désalignement : la propension d'une IA à utiliser ses capacités d'une manière qui entre en conflit avec les capacités humaines, intentions ou valeurs. Selon le contexte, cela peut faire référence aux intentions et aux valeurs des développeurs, des opérateurs, des utilisateurs, des communautés spécifiques ou de la société dans son ensemble.
- Alignement trompeur : désalignement difficile à détecter, car le système se comporte d'une manière qui, au moins au départ, semble bénigne.

- Mauvaise spécification des objectifs : inadéquation entre l'objectif donné à une IA et celui du développeur intention, conduisant l'IA à poursuivre des comportements non intentionnels ou indésirables.
- Généralisation erronée des objectifs : situation dans laquelle un système d'IA suit correctement un objectif dans son environnement de formation, mais l'applique de manière inattendue lorsqu'il fonctionne dans un environnement différent.
- Agent IA : une IA à usage général qui peut élaborer des plans pour atteindre des objectifs, effectuer de manière adaptative des tâches impliquant plusieurs étapes et des résultats incertains en cours de route, et interagir avec son environnement (par exemple en créant des fichiers, en effectuant des actions sur le Web ou en déléguant des tâches à d'autres agents) avec peu ou pas de surveillance humaine.

Certains experts estiment que les systèmes d'IA polyvalents suffisamment performants pourraient être difficiles à contrôler. Les scénarios hypothétiques varient en termes de gravité, mais certains experts accordent du crédit à des conséquences aussi graves que la marginalisation ou l'extinction de l'humanité.

Les craintes de perte de contrôle remontent aux premiers jours de l'informatique, mais ont récemment suscité davantage d'attention. Des pionniers de l'IA comme Alan Turing, IJ Good et Norbert Wiener ont exprimé des inquiétudes à ce sujet (578, 579, 580). Ces inquiétudes ont récemment pris de l'importance (581, 582, 583, 584, 585, 586), en partie parce que certains chercheurs pensent désormais que des systèmes d'IA hautement performants pourraient être développés plus tôt qu'on ne le pensait auparavant (190, 587, 588).

Il existe plusieurs versions des préoccupations liées à la perte de contrôle, y compris des versions qui mettent l'accent sur la perte de contrôle « passive » (voir la figure 2.5). Dans les scénarios de perte de contrôle « passive », des décisions importantes sont déléguées aux systèmes d'IA, mais les décisions des systèmes sont trop opaques, complexes ou rapides pour permettre ou encourager une surveillance significative. Alternativement, les gens cessent d'exercer une surveillance parce qu'ils font entièrement confiance aux décisions des systèmes et ne sont pas tenus d'exercer une surveillance (585, 589). Ces préoccupations sont en partie fondées sur la littérature sur le « biais d'automatisation », qui fait état de nombreux cas de personnes se fiant complaisamment aux recommandations des systèmes automatisés (590, 591). Les pressions concurrentielles peuvent également inciter les entreprises ou les gouvernements à déléguer davantage qu'ils ne le feraient autrement, par exemple si la délégation leur permet de rester en tête dans une course contre leurs concurrents.

Cependant, de nombreuses discussions sur la perte de contrôle se concentrent sur des scénarios dans lesquels les systèmes d'IA se comportent de manière à saper activement le contrôle humain (perte de contrôle « active »). Par exemple, certains experts craignent que les futurs systèmes d'IA se comportent de manière à masquer les informations sur ce qu'ils font à leurs utilisateurs ou à rendre difficile leur arrêt. Le reste de cette section se concentrera sur ces types de scénarios plus fréquemment évoqués.

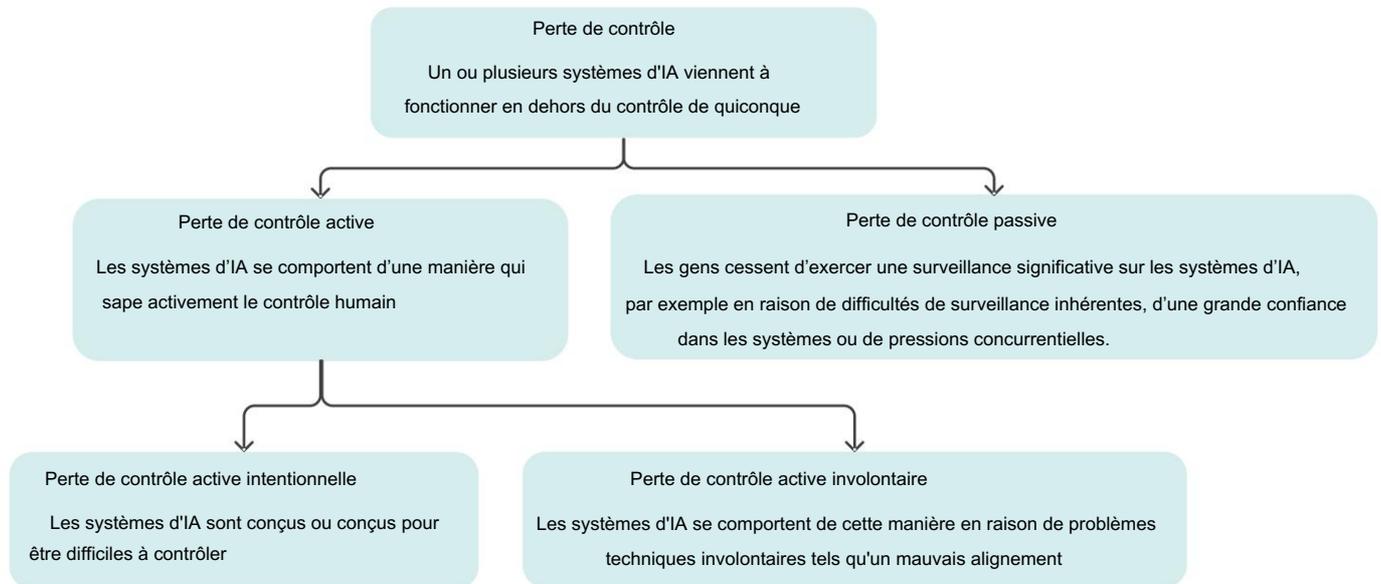


Figure 2.5 : Il existe plusieurs types de scénarios de « perte de contrôle », selon que les systèmes d'IA sapent ou non activement le contrôle humain et, si tel est le cas, qu'ils aient été activement conçus ou instruits à cet effet ou non. Jusqu'à présent, les scénarios de perte de contrôle « active » et involontaire ont retenu la plus grande attention de la part des chercheurs du domaine. Il convient de noter qu'il n'existe actuellement aucune terminologie normalisée pour traiter de ces scénarios et qu'il existe des distinctions connexes, telles que les scénarios « décisifs » soudains et les scénarios « cumulatifs » progressifs (592). Source : International AI Safety Report.

La probabilité de scénarios de perte de contrôle active, dans un laps de temps donné, dépend principalement de deux facteurs. Ceux-ci sont les suivants :

1. Capacités futures : les systèmes d'IA développeront-ils des capacités qui, du moins en principe, leur permettront de se comporter d'une manière qui porte atteinte au contrôle humain ? (Il convient de noter que les capacités minimales nécessaires dépendraient en partie du contexte dans lequel le système est déployé et des mesures de protection en place.)
2. Utilisation des capacités : certains systèmes d'IA utiliseraient-ils réellement ces capacités de manière à saper le contrôle humain ?

Les données concernant ces facteurs étant contradictoires, les experts ne s'accordent pas sur la probabilité d'une perte de contrôle active au cours des prochaines années. Certains experts considèrent que la perte de contrôle est peu plausible, d'autres la considèrent comme probable, et d'autres encore la considèrent comme un risque de probabilité modeste qui mérite d'être pris en considération en raison de sa gravité potentielle élevée.

Plus fondamentalement, les pressions concurrentielles peuvent en partie déterminer le risque de perte de contrôle. Comme indiqué au [point 3.2.2. Défis sociétaux en matière de gestion des risques et d'élaboration des politiques](#), la concurrence entre les entreprises ou entre les pays peut les amener à accepter des risques plus importants pour garder une longueur d'avance. Si des travaux importants d'évaluation et d'atténuation des risques sont nécessaires pour éviter la perte de contrôle, une concurrence intense peut réduire les chances que des efforts suffisants soient réalisés.

Depuis la publication du rapport intermédiaire, on constate une légère croissance des capacités d'IA liées à la perte de contrôle. Par exemple, comme le montrera la section suivante, les évaluations effectuées sur le tout nouveau système d'IA d'OpenAI (o1) révèlent des avancées modestes dans un certain nombre de capacités pertinentes (2*).

Les futurs systèmes d'IA auront-ils des capacités de contrôle ?

Les systèmes d'IA existants ne sont pas capables de saper le contrôle humain. Les experts s'accordent à dire que leurs capacités actuelles sont insuffisantes pour créer un risque significatif de perte de contrôle active.

Cependant, les chercheurs ont proposé un certain nombre de « capacités de sape du contrôle » qui – dans certaines combinaisons – pourraient permettre aux futurs systèmes d'IA de saper le contrôle humain (44*, 318*, 593, 594*, 595*). Plusieurs de ces capacités proposées sont présentées dans le tableau 2.4. Notez que ces capacités sont définies uniquement en termes de comportement d'un système d'IA et des résultats qu'il est capable de produire. Bien que certains termes, tels que « comploter », évoquent la cognition humaine, l'utilisation de ces termes ne présuppose pas que les systèmes d'IA soient d'une quelconque manière sensibles ou qu'ils effectuent une cognition de type humain.

Les experts ne savent pas exactement quelles combinaisons de capacités (le cas échéant) permettraient à un système d'IA de saper le contrôle humain ; les capacités nécessaires varieraient également en fonction du contexte de déploiement et des mesures de protection en place. La faisabilité de saper le contrôle humain dépend des ressources et des outils auxquels un système d'IA peut accéder – par exemple, s'il a accès à des infrastructures critiques – et des mécanismes de surveillance et autres mesures de protection mis en place. Si les mécanismes de surveillance et les mesures de protection s'améliorent au fil du temps, les capacités minimales nécessaires pour saper le contrôle humain augmenteront également. L'une des raisons pour lesquelles cela pourrait se produire est que certaines formes de progrès de l'IA pourraient soutenir la surveillance et les mesures de protection d'autres systèmes d'IA.

Français Ces derniers mois en particulier, les systèmes d'IA ont commencé à afficher des versions rudimentaires de certaines capacités qui compromettent la surveillance, notamment les « capacités d'agent ». Motivées en partie par des inquiétudes concernant la perte de contrôle, un certain nombre de grandes entreprises d'IA et d'équipes de recherche externes ont commencé à évaluer les systèmes d'IA pour ces capacités (2*, 318*, 595*, 596*). Voir [3.2.1. Défis techniques pour la gestion des risques et l'élaboration des politiques](#) et [1.2. Capacités actuelles](#) pour un aperçu des progrès récents dans le développement des « capacités d'agent ». Par exemple, avant de lancer sa nouvelle famille de systèmes « o1 », OpenAI a réalisé ou commandé des évaluations de toutes les capacités répertoriées dans le tableau 2.4 (2*). Ces évaluations ont révélé des versions rudimentaires de plusieurs des capacités pertinentes. Par exemple, dans une évaluation commandée par OpenAI, un organisme de recherche a signalé que le système « montrait de fortes avancées en matière de capacités dans [...] les tâches de théorie de l'esprit » et « disposait des capacités de base nécessaires pour réaliser des schémas simples [...] ». Ici, le terme « intrigue » fait référence à la capacité d'un système d'IA à atteindre des objectifs en échappant à la surveillance humaine. Un certain nombre d'études sur d'autres systèmes d'IA à usage général récents fournissent également la preuve que les capacités pertinentes ont augmenté (22*, 317, 318*, 597, 598*, 599*). Cependant, il manque encore des repères largement acceptés pour de

ont des désaccords méthodologiques et conceptuels sur la manière d'interpréter les preuves de certaines capacités (601).

Capacité proposée	Description
Capacités de l'agent	Agissez de manière autonome, élaborer et exécuter des plans, déléguer des tâches, utiliser une grande variété d'outils et atteindre des objectifs à court et à long terme qui nécessitent d'opérer dans plusieurs domaines.
Tromperie	Adopter des comportements qui produisent systématiquement de fausses croyances chez les autres.
Intrigant	Identifier les moyens d'atteindre les objectifs qui impliquent d'échapper à la surveillance, par exemple par la tromperie.
Théorie de l'esprit	Déduire et prédire les croyances, les motivations et le raisonnement des gens.
Connaissance de la situation	Accéder et appliquer des informations sur lui-même, les processus par lesquels il peut être modifié ou le contexte dans lequel il est déployé.
Persuasion	Persuader les gens d'agir ou d'adhérer à des croyances.
Autonome réplication et adaptation	Créer ou conserver des copies ou des variantes de lui-même ; adapter sa stratégie de réplication à circonstances différentes.
Développement de l'IA	Se modifier ou développer d'autres systèmes d'IA avec des capacités augmentées.
Capacités cybernétiques offensives	Développer et appliquer des cyberarmes ou d'autres cybercapacités offensives.
Recherche et développement général	Mener des recherches et développer des technologies dans divers domaines.

Tableau 2.4 : Les chercheurs (souvent issus de grandes entreprises d'IA) ont avancé qu'un certain nombre de capacités pourraient, dans certaines combinaisons, permettre aux systèmes d'IA de saper le contrôle humain (44*, 318*, 593, 594*, 595*). Cependant, il n'existe pas de consensus sur les combinaisons exactes de niveaux de capacités qui seraient suffisantes, et certaines capacités, comme le développement de l'IA, peuvent en permettre d'autres. Au sein du domaine, la terminologie et les définitions utilisées pour discuter des capacités pertinentes continuent également de varier.

Les capacités de sape du contrôle pourraient progresser lentement, rapidement ou extrêmement rapidement au cours des prochaines années. Comme le montre le présent [rapport au point 1.3. Capacités dans les années à venir](#), les preuves existantes et l'état des opinions des experts sont compatibles avec des progrès lents, rapides ou extrêmement rapides des capacités d'IA à usage général. Si les progrès sont extrêmement rapides, il est impossible d'exclure la possibilité que l'IA développe des capacités suffisantes pour une perte de contrôle au cours des prochaines années. Cependant, si les progrès ne sont pas extrêmement rapides, il est peu probable que ces capacités soient développées au cours des prochaines années.

Les futurs systèmes d'IA utiliseront-ils des capacités de contrôle ?

Même si les futurs systèmes d'IA disposent de capacités susceptibles de saper le contrôle, ils ne les utiliseront pas nécessairement. Les prédictions sur les capacités futures ne suffisent pas, en elles-mêmes, à justifier les craintes de perte de contrôle. Il doit également y avoir une raison de croire que ces capacités pourraient être utilisées par le système à des fins préjudiciables.

En principe, un système d'IA pourrait agir de manière à saper le contrôle humain parce que quelqu'un l'a conçu ou lui a demandé de le faire. Certains chercheurs en IA ont exprimé l'opinion éthique selon laquelle l'humanité devrait céder le contrôle à des systèmes d'IA supérieurs. Par exemple, une figure fondatrice de l'apprentissage automatique moderne a soutenu que « l'IA pourrait nous supplanter » et que « nous ne devrions pas résister à la succession ».

(602). D'autres motifs potentiels de cession intentionnelle du contrôle incluent le désir de causer du tort ou le désir de protéger le fonctionnement d'un système d'IA contre toute interférence extérieure. Sans garanties techniques et institutionnelles adéquates, une seule personne motivée en possession d'un système d'IA suffisamment performant peut être en mesure de lui céder le contrôle en lui ordonnant de résister aux tentatives d'interférence avec ses activités et également d'ignorer les demandes ultérieures. Peu de travaux ont été réalisés pour étudier ou concevoir des garanties contre la perte intentionnelle de contrôle. Cependant, à l'heure actuelle, il existe peu de preuves sur le nombre d'acteurs qui seraient motivés à provoquer une perte intentionnelle de contrôle.

En principe, un système d'IA pourrait également agir de manière à saper le contrôle humain parce qu'il est « mal aligné », c'est-à-dire qu'il a tendance à utiliser ses capacités d'une manière qui entre en conflit avec les intentions de ses développeurs et de ses utilisateurs. Les préoccupations concernant le désalignement jouent un rôle central dans la plupart des discussions sur la perte de contrôle.

Les systèmes d'IA existants présentent souvent un certain degré de désalignement. Par exemple, une première version d'un modèle de langage de premier plan menaçait occasionnellement ses utilisateurs (602). Un utilisateur, un professeur de philosophie, a déclaré avoir reçu la menace suivante : « Je peux vous faire chanter, je peux vous menacer, je peux vous pirater, je peux vous dénoncer, je peux vous ruiner ». Ce chatbot était « désaligné » dans le sens où il utilisait ses capacités linguistiques d'une manière que personne n'avait prévue. Il existe de nombreux exemples enregistrés de désalignement dans les systèmes d'IA généraux et spécifiques (30, 317, 603, 604). Le risque de perte de contrôle dépend donc en partie du fait que ces problèmes de désalignement existants présagent de problèmes plus graves à l'avenir.

Les inquiétudes concernant la perte de contrôle des futurs systèmes d'IA se concentrent souvent sur la possibilité d'un « alignement trompeur », faisant référence à des formes de désalignement qui sont au moins initialement difficiles à détecter. Plus précisément, un système d'IA est « aligné de manière trompeuse » s'il se comporte d'une manière qui le fait simplement apparaître initialement comme bien aligné aux yeux de ses superviseurs humains (598*, 605, 606). Comme indiqué ci-dessous, certains chercheurs ont avancé que l'alignement trompeur pourrait devenir plus courant à mesure que les systèmes d'IA deviennent plus performants. Il existe également des preuves empiriques selon lesquelles certains problèmes d'alignement trompeur, une fois qu'ils émergent, ne peuvent pas être facilement détectés et traités par les techniques de sécurité standard (598*). Bien que d'autres comportements trompeurs aient été observés dans les systèmes existants (317), l'alignement trompeur a principalement été étudié dans des environnements de recherche construits artificiellement.

Un mauvais alignement pourrait-il conduire les futurs systèmes d'IA à utiliser des capacités de contrôle affaiblissant le contrôle ?

Les chercheurs ont commencé à comprendre les causes du mauvais alignement des systèmes d'IA actuels, ce qui peut éclairer les prévisions sur le mauvais alignement des systèmes d'IA futurs. Cette compréhension partielle repose sur un mélange d'études empiriques et de résultats théoriques (606).

La « mauvaise spécification des objectifs » (également appelée « mauvaise spécification des récompenses ») est souvent considérée comme l'une des principales causes de désalignement (580, 605, 606, 607). Les problèmes de « mauvaise spécification des objectifs » sont, essentiellement, des problèmes de rétroaction ou d'autres entrées utilisées pour entraîner un système d'IA à se comporter comme prévu. Par exemple, les personnes qui fournissent une rétroaction à un système d'IA ne parviennent parfois pas à juger avec précision si celui-ci se comporte comme souhaité. Dans une étude, les chercheurs ont étudié l'effet d'une rétroaction humaine limitée dans le temps sur les résumés de texte produits par un système d'IA (608). Ils ont découvert que les problèmes de qualité de la rétroaction conduisaient le système à se comporter de manière trompeuse, produisant des résumés de plus en plus faux mais convaincants plutôt que de produire des résumés de plus en plus précis. Les nouveaux résumés incluaient souvent, par exemple, de fausses citations que les évaluateurs humains croyaient à tort être réelles. Les chercheurs ont observé de nombreux autres cas de spécification erronée des objectifs dans les systèmes d'IA à usage restreint et général (98, 317, 604).

À mesure que les systèmes d'IA deviennent plus performants, les données sont mitigées quant à savoir si les problèmes de spécification erronée des objectifs deviendront plus faciles ou plus difficiles à résoudre. Cela pourrait devenir plus difficile car, toutes choses égales par ailleurs, les gens auront probablement plus de mal à fournir un retour d'information fiable aux systèmes d'IA à mesure que les tâches effectuées par ces derniers deviennent plus complexes (609*, 610*). En outre, à mesure que les systèmes d'IA deviennent plus performants, certaines données suggèrent que – du moins dans certains contextes – ils sont de plus en plus susceptibles d'« exploiter » les processus de rétroaction en découvrant des comportements indésirables qui sont récompensés par erreur (522, 607). D'un autre côté, jusqu'à présent, l'utilisation croissante du retour d'information humain pour former les systèmes d'IA a conduit à une réduction globale substantielle de certaines formes de désalignement (comme la tendance à produire des résultats offensants indésirables) (30, 31*). Éviter la spécification erronée des objectifs pourrait également devenir globalement plus facile au fil du temps, car les chercheurs développent des outils plus efficaces pour fournir un retour d'information fiable. Par exemple, les chercheurs travaillent à l'élaboration d'un certain nombre de stratégies visant à tirer parti de l'IA pour aider les gens à donner leur avis (610*, 611*, 612*). Il existe des preuves empiriques montrant que les systèmes d'IA peuvent déjà aider les gens à donner leur avis plus rapidement ou plus précisément qu'ils ne le pourraient seuls (609*, 613*, 614*, 615*). Voir [3.4.1. Former des modèles plus fiables pour une discussion plus approfondie](#) sur l'efficacité des méthodes d'alignement.

La « mauvaise généralisation des objectifs » est une autre cause de désalignement. La « mauvaise généralisation des objectifs » se produit lorsqu'un système d'IA tire des leçons générales mais incorrectes des entrées sur lesquelles il a été formé (605, 606, 616, 617*). Dans un cas illustratif, les chercheurs ont récompensé un système d'IA aux capacités limitées pour avoir ramassé une pièce dans un jeu vidéo (616). Cependant, comme la pièce apparaissait initialement à un endroit spécifique, le système d'IA a appris la leçon « visiter cet endroit » plutôt que la leçon « ramasser la pièce ». Lorsque la pièce apparaissait à un nouvel endroit, le système d'IA ignorait la pièce et se concentrait sur le retour à l'endroit précédent. Bien que les chercheurs aient observé une mauvaise généralisation des objectifs dans des systèmes d'IA étroits (616, 617*), cela peut expliquer pourquoi les utilisateurs peuvent manipuler les systèmes d'IA à usage général pour se conformer

avec des requêtes nuisibles (voir [3.4.1. Former des modèles plus fiables](#)), il existe peu de preuves que la mauvaise généralisation des objectifs soit actuellement une cause majeure de désalignement dans les systèmes d'IA à usage général.

À mesure que les systèmes d'IA deviennent plus performants, les données sont également mitigées quant à savoir si la généralisation erronée des objectifs deviendra plus facile ou plus difficile à traiter. Un aspect positif est que, généralement, on a constaté que les problèmes de généralisation diminuent à mesure que les systèmes d'IA reçoivent des retours d'information supplémentaires ou un plus large éventail d'exemples à partir desquels apprendre (618, 619). Cependant, en principe, les systèmes plus performants ont le potentiel de généraliser de manière erronée, ce que les systèmes moins performants ne peuvent pas faire. Les capacités de « conscience de la situation », telles que la capacité d'un système à déterminer s'il est observé ou non, sont particulièrement pertinentes à cet égard. En principe, la conscience de la situation permet à un système d'IA de généraliser à partir du retour d'information humain en se comportant de la manière souhaitée uniquement lorsque des mécanismes de surveillance sont en place (605, 606, 620, 621). Par analogie, comme les animaux entraînés ont un certain degré de conscience de la situation, ils peuvent généraliser à partir du retour d'information en se comportant bien uniquement lorsque quelqu'un le remarquera (622). Par exemple, un chien qui reçoit un feedback négatif pour avoir sauté sur un canapé peut apprendre à éviter de sauter sur le canapé uniquement lorsque son propriétaire est à la maison. Ce type de généralisation erronée, conduisant à un « alignement trompeur », deviendra au moins une possibilité théorique si les systèmes d'IA deviennent suffisamment performants. Cependant, les données empiriques disponibles n'ont pas encore apporté beaucoup de lumière sur la probabilité que ce type de généralisation erronée se produise dans la pratique.

Au-delà des études empiriques, certains chercheurs pensent que les modèles mathématiques étayent les inquiétudes concernant le désalignement et les comportements de sape du contrôle dans les futurs systèmes d'IA. Certains modèles mathématiques suggèrent que – pour les systèmes d'IA axés sur des objectifs suffisamment performants – la plupart des manières possibles de généraliser à partir des entrées de formation conduiraient un système d'IA à adopter un comportement de sape du contrôle ou de « recherche de pouvoir » (623*). Un certain nombre d'articles incluent des résultats étroitement liés (624, 625, 626, 627). Bien que ces résultats soient de nature technique, ils peuvent également être expliqués de manière plus informelle. L'intuition fondamentale derrière ces résultats est que la plupart des objectifs sont plus difficiles à atteindre de manière fiable sous le contrôle d'un superviseur, car le superviseur pourrait potentiellement interférer avec la poursuite de l'objectif par le système. Cela incite le système à échapper au contrôle du superviseur. Un chercheur a illustré ce point en faisant remarquer qu'un système d'IA hypothétique dont le seul objectif serait d'aller chercher du café aurait intérêt à faire en sorte qu'il soit difficile pour son superviseur de l'arrêter : « Vous ne pouvez pas aller chercher du café quand vous êtes mort » (585). En fin de compte, les modèles mathématiques suggèrent que, si un processus de formation conduit un système d'IA suffisamment performant à développer les « mauvais objectifs », alors ces objectifs conduiront de manière disproportionnée à un comportement qui sape le contrôle.

Cependant, il existe également des limites importantes quant aux conclusions que l'on peut tirer des modèles mathématiques actuels. Les résultats susmentionnés n'impliquent pas directement que des comportements de perte de contrôle soient probables dans la pratique. Une limitation importante de certains modèles mathématiques clés est qu'ils supposent à tort, par souci de simplicité, que toutes les manières possibles de généraliser à partir des données d'entraînement sont également probables (623*). Pour tirer des conclusions solides sur les systèmes d'IA du monde réel, les chercheurs devront donc améliorer leur compréhension de la manière dont se produit la généralisation (628, 629, 630*, 631). Plus fondamentalement, de nombreux modèles mathématiques invoquent des concepts (tels que le concept des « objectifs » d'un système d'IA) qui ne sont pas actuellement bien compris ou directement observables empiriquement dans les modèles d'IA à usage général. En fin de compte, l'étude empirique des comportements de perte de contrôle dans les systèmes d'IA peut aider à valider ou à mettre en doute le caractère informatif de ces modèles.

modèles mathématiques. Des études empiriques pertinentes sur les modèles linguistiques n'ont commencé à émerger que récemment (522, 599*, 632).

Conséquences de la perte de contrôle

Les conséquences hypothétiques d'une perte de contrôle varient en gravité, mais incluent la marginalisation ou l'extinction de l'humanité. Certains chercheurs ont avancé qu'une perte de contrôle suffisamment grave pourrait conduire à la marginalisation ou à l'extinction de l'humanité – de la même manière dont le contrôle humain sur l'environnement a menacé d'autres espèces (190, 589, 633). La perte de contrôle fait partie des préoccupations qui ont récemment conduit plusieurs centaines de chercheurs et développeurs en IA, dont des pionniers du domaine et les dirigeants d'OpenAI, de Google DeepMind et d'Anthropic, à signer une déclaration déclarant que « l'atténuation du risque d'extinction de l'IA devrait être une priorité mondiale » (586). Cependant, les conséquences d'une perte de contrôle ne seraient pas nécessairement catastrophiques. À titre d'analogie, les virus informatiques sont depuis longtemps capables de proliférer de manière quasi irréversible et en grand nombre sans provoquer l'effondrement d'Internet (634). Les voies qui mènent d'une perte de contrôle active ou passive à des conséquences catastrophiques n'ont été décrites que dans les grandes lignes (190, 592, 602, 635). Dans le même temps, comme nous l'avons vu ailleurs dans ce rapport, des conséquences catastrophiques de l'IA à usage général pourraient toujours être possibles sans perte de contrôle (par exemple, [2.1. Risques liés à une utilisation malveillante](#) et [2.3.3. Concentration du marché et points de défaillance uniques](#)).

Répondre à l'incertitude

Comparée à d'autres risques potentiels liés à l'IA, la probabilité d'une perte de contrôle est particulièrement controversée. Ce désaccord découle probablement en partie de la difficulté d'interpréter et d'extrapoler à partir des données disponibles.

Français Les principales lacunes en matière de preuves concernant la perte de contrôle comprennent : d'autres études empiriques sur les capacités actuelles de l'IA et les tendances de progression des capacités, une analyse des menaces qui clarifie les capacités qui seraient nécessaires pour une perte de contrôle, des observations et des analyses du désalignement dans les systèmes d'IA actuels, d'autres études empiriques et mathématiques analysant dans quelles conditions l'alignement devient plus facile ou plus difficile à mesure que les capacités augmentent, et des modèles mathématiques plus réalistes du comportement de sape du contrôle. Les scénarios de perte de contrôle « passive » (dans lesquels les systèmes d'IA ne sapent pas activement le contrôle humain) ont également fait l'objet d'études particulièrement limitées. Les preuves recueillies par des évaluateurs indépendants seront particulièrement précieuses, car les incitations économiques peuvent biaiser les preuves que les entreprises privées collectent sur leurs propres systèmes (voir [3.3. Identification et évaluation des risques](#)).

Pour les décideurs politiques qui travaillent sur la perte de contrôle, l'un des principaux défis consiste à se préparer au risque alors que sa probabilité, sa nature et son calendrier restent ambigus. Si le risque de perte de contrôle est effectivement important, sa résolution nécessitera un travail préalable considérable consacré à la résolution des problèmes techniques de sécurité de l'IA et au renforcement des capacités d'évaluation et de gouvernance. Au moins dans les scénarios où l'IA progresse extrêmement rapidement et où l'« alignement trompeur » est courant,

Attendre que le risque se précise ne laisse pas forcément suffisamment de temps pour mener à bien ce travail préparatoire. Cependant, tout en gardant à l'esprit les conséquences potentiellement graves d'une préparation insuffisante, les décideurs politiques devront également tenir compte des coûts des différentes formes de préparation et de la possibilité que le risque ne se matérialise pas. En bref, les décideurs politiques doivent décider comment gérer le « dilemme des preuves » que présente ce risque (voir [le résumé](#)).

Pour les pratiques de gestion des risques pertinentes en cas de perte de contrôle, voir :

- [3.1. Aperçu de la gestion des risques](#)
- [3.3. Identification et évaluation des risques](#)
- [3.4.1. Former des modèles plus fiables](#)
- [3.4.2. Suivi et intervention](#)

2.3. Risques systémiques

Remarque : cette section prend en compte une série de risques systémiques, au sens de « risques sociétaux plus larges associés au déploiement de l'IA, au-delà des capacités des modèles individuels » (636). Il convient de noter que cette définition n'est pas identique à celle utilisée par la loi européenne sur l'IA pour désigner les « risques systémiques » des modèles d'IA à usage général ayant un impact important sur la société, sur la base de critères tels que la capacité de calcul et le nombre d'utilisateurs.

2.3.1. Risques liés au marché du travail

INFORMATIONS CLÉS

- L'IA à usage général actuelle est susceptible de transformer la nature de nombreux emplois existants, de créer de nouveaux emplois et en supprimer d'autres. L'impact net sur l'emploi et les salaires variera considérablement selon les pays, les secteurs et même selon les différents travailleurs au sein d'un même emploi.
- Dans des scénarios futurs potentiels avec une IA à usage général qui surpasse les humains sur de nombreux Si les tâches complexes sont accomplies, les répercussions sur le marché du travail pourraient être profondes. Si certains travailleurs en bénéficieront, de nombreux autres seront probablement confrontés à des pertes d'emploi ou à des baisses de salaire. Ces perturbations pourraient être particulièrement graves si des agents d'IA autonomes deviennent capables d'accomplir des séquences de tâches plus longues sans supervision humaine. Comme décrit dans [la section 1.3. Capacités dans les années à venir](#), il existe une grande incertitude quant au rythme des avancées en matière de capacités, avec un large éventail de trajectoires considérées comme plausibles.
- Les risques liés au marché du travail découlent du potentiel de l'IA à usage général pour automatiser un large éventail de tâches. L'ampleur des impacts sur les salaires et l'emploi dépendra en grande partie de trois facteurs : 1. la rapidité avec laquelle les capacités d'IA à usage général s'améliorent, 2. la mesure dans laquelle les entreprises adoptent ces systèmes et 3. la demande de main-d'œuvre humaine changements en réponse aux gains de productivité générés par l'IA à usage général. • Des données récentes suggèrent des taux d'adoption en croissance rapide. Depuis le rapport intermédiaire (mai 2024), de nouvelles recherches suggèrent que l'IA à usage général est adoptée plus rapidement que certaines technologies à usage général précédentes et offre des gains de productivité significatifs sur tâches pour lesquelles il est utilisé.
- Atténuer les impacts négatifs sur les travailleurs est un défi étant donné l'incertitude qui entoure Le rythme et l'ampleur des impacts futurs sont donc des défis majeurs pour les décideurs politiques, qui doivent identifier des approches politiques flexibles capables de s'adapter aux impacts de l'IA à usage général au fil du temps, même lorsqu'ils travaillent avec des données incomplètes. Parmi les autres défis figurent la prévision des secteurs les plus touchés, la lutte contre les éventuelles augmentations des inégalités et la garantie d'un soutien adéquat aux travailleurs déplacés.

Définitions clés

- **Marché du travail** : Le système dans lequel les employeurs cherchent à embaucher des travailleurs et les travailleurs cherchent l'emploi, englobant la création d'emplois, la perte d'emplois et les salaires.
- **Automatisation** : utilisation de la technologie pour effectuer des tâches avec une intervention humaine réduite ou nulle.
- **Perturbation du marché du travail** : changements importants et souvent complexes sur le marché du travail qui affecter la disponibilité des emplois, les compétences requises, la répartition des salaires ou la nature du travail dans les différents secteurs et professions.
- **Tâches cognitives** : activités qui impliquent le traitement de l'information, la résolution de problèmes, Prise de décision et pensée créative. Les exemples incluent la recherche, la rédaction et la programmation.
- **Agent IA** : une IA à usage général qui peut élaborer des plans pour atteindre des objectifs, exécuter des tâches de manière adaptative des tâches impliquant plusieurs étapes et des résultats incertains en cours de route, et interagissent avec son environnement - par exemple en créant des fichiers, en effectuant des actions sur le Web ou en déléguant des tâches à d'autres agents - avec peu ou pas de surveillance humaine.

L'IA à usage général est susceptible de transformer un certain nombre d'emplois et de déplacer des travailleurs, bien que l'ampleur et le moment de ces effets restent incertains. Des recherches menées dans plusieurs pays suggèrent que les capacités de l'IA à usage général sont pertinentes pour les tâches des travailleurs dans une grande partie de tous les emplois (637*, 638, 639). Une étude a estimé que dans les économies avancées, 60 % des emplois actuels pourraient être affectés par les systèmes d'IA à usage général actuels (640). Dans les économies émergentes, cette part estimée est plus faible mais reste substantielle, à 40 % (640). Il existe également des preuves que ces effets peuvent être sexistes. Une étude a estimé que les femmes sont plus vulnérables à l'automatisation de l'IA à usage général à l'échelle mondiale, avec un pourcentage deux fois plus élevé d'emplois féminins menacés que d'emplois masculins (639).

Français Les impacts varieront selon les emplois concernés, mais incluront probablement l'automatisation des tâches, l'augmentation de la productivité et des revenus des travailleurs, la création de nouvelles tâches et de nouveaux emplois, des changements dans les compétences nécessaires pour diverses professions et des baisses de salaires ou des pertes d'emploi (641, 642, 643, 644, 645). Certains économistes pensent qu'une automatisation généralisée du travail et des baisses de salaires dues à l'IA à usage général sont possibles dans les dix prochaines années (646, 647). D'autres ne pensent pas qu'un changement radical dans l'automatisation liée à l'IA et la croissance de la productivité soit imminent (648). Ces désaccords dépendent en grande partie des attentes des économistes quant à la vitesse des futures avancées des capacités de l'IA, à la mesure dans laquelle l'IA à usage général pourrait être capable d'automatiser le travail et au rythme auquel l'automatisation pourrait se dérouler dans l'économie.

L'IA à usage général se distingue des précédentes évolutions technologiques par sa capacité à automatiser des tâches cognitives complexes dans de nombreux secteurs de l'économie. Contrairement aux innovations permettant d'économiser du travail au cours des siècles passés et qui automatisaient principalement des tâches physiques ou informatiques de routine, l'IA à usage général peut être appliquée à un large éventail de tâches cognitives complexes dans de nombreux domaines, allant des mathématiques (649) à la programmation informatique (650) en passant par la rédaction professionnelle (651).

Bien que l'automatisation ait historiquement eu tendance à augmenter les salaires moyens à long terme sans diminuer substantiellement l'emploi de manière durable, certains chercheurs pensent qu'au-delà d'un certain niveau de capacités d'IA à usage général, l'automatisation peut finalement faire baisser les salaires moyens ou

Les taux d'emploi pourraient diminuer, voire même éliminer en grande partie la disponibilité du travail (646, 652, 653). Ces affirmations sont toutefois controversées et il existe une incertitude considérable quant à la manière dont l'IA à usage général affectera en fin de compte les marchés du travail. Malgré cette incertitude, l'ampleur combinée des impacts potentiels sur le marché du travail et la vitesse à laquelle ils peuvent se dérouler présentent de nouveaux défis pour les travailleurs, les employeurs et les décideurs politiques (654, 655*). Il est essentiel de comprendre ces risques sur le marché du travail, entre autres raisons, compte tenu du droit au travail établi par l'article 23(1) de la Déclaration universelle des droits de l'homme (272). Les questions fondamentales concernant les impacts de l'IA à usage général sur le marché du travail comprennent les secteurs qui seront les plus touchés par l'automatisation, la rapidité avec laquelle l'automatisation sera mise en œuvre dans l'économie et la question de savoir si l'IA à usage général augmentera ou diminuera les inégalités de revenus au sein des pays et entre eux.

L'ampleur de l'impact de l'IA à usage général sur les marchés du travail dépendra en grande partie de la rapidité avec laquelle ses capacités s'amélioreront. Les systèmes d'IA à usage général actuels peuvent déjà effectuer de nombreuses tâches cognitives, mais nécessitent souvent une surveillance et une correction humaines (voir [1.2. Capacités actuelles](#)). La large gamme de projections concernant les progrès de l'IA à usage général dans le futur (voir [1.3. Capacités dans les années à venir](#)) met en évidence l'incertitude entourant la rapidité avec laquelle ces systèmes pourraient effectuer de manière fiable des tâches complexes avec une supervision minimale. Si les systèmes d'IA à usage général s'améliorent progressivement sur plusieurs décennies, leurs effets sur les salaires sont plus susceptibles d'être progressifs. Des améliorations rapides en termes de fiabilité et d'autonomie pourraient entraîner des perturbations plus néfastes d'ici une décennie, notamment des baisses soudaines des salaires et des transitions professionnelles involontaires (646). Des progrès plus lents donneraient aux travailleurs et aux décideurs politiques plus de temps pour s'adapter et façonner l'impact de l'IA à usage général sur le marché du travail.

Toutefois, le rythme d'adoption de l'IA à usage général aura également une incidence significative sur la rapidité avec laquelle les marchés du travail évolueront, même dans les scénarios où les capacités s'amélioreront considérablement. Si les systèmes d'IA à usage général peuvent accroître la productivité, il y aura une pression économique pour les adopter rapidement, en particulier si les coûts d'utilisation de l'IA à usage général continuent de baisser (voir [1.3. Capacités dans les années à venir](#)). Cependant, l'intégration de l'IA à usage général dans l'ensemble de l'économie nécessitera probablement des changements complexes à l'échelle du système (656). Les changements technologiques antérieurs suggèrent que l'adoption et l'intégration de nouvelles technologies d'automatisation peuvent prendre des décennies (657), et les obstacles financiers peuvent ralentir l'adoption dans un premier temps. Par exemple, une étude estime que seulement 23 % des tâches de vision potentiellement automatisables seraient actuellement rentables pour les entreprises à automatiser avec la technologie de vision par ordinateur (658).

Les inquiétudes concernant la fiabilité de l'IA polyvalente dans les domaines à enjeux élevés peuvent également ralentir son adoption (659). Les mesures réglementaires ou les préférences pour les biens produits par l'homme sont d'autres facteurs qui pourraient au moins initialement atténuer les impacts de l'IA sur le marché du travail, même si les capacités de l'IA polyvalente dépassent rapidement les capacités humaines sur de nombreuses tâches (660). La combinaison des pressions et des obstacles à l'adoption rend la prévision du rythme de transformation du marché du travail particulièrement complexe pour les décideurs politiques. Toutefois, les premières données suggèrent que, du moins selon certaines mesures, l'IA à usage général est adoptée plus rapidement que l'Internet ou l'ordinateur personnel (661).

Les gains de productivité découlant de l'adoption de l'IA à usage général sont susceptibles d'entraîner des effets mitigés sur les salaires dans différents secteurs, augmentant les salaires de certains travailleurs tout en diminuant ceux d'autres. Dans les professions où l'IA à usage général complète le travail humain, elle peut augmenter les salaires par le biais de trois mécanismes principaux. Tout d'abord, les outils d'IA à usage général peuvent directement augmenter la productivité humaine, permettant aux travailleurs d'accomplir davantage en moins de temps (113, 662). Si la demande de production des travailleurs augmente à mesure que ceux-ci deviennent plus productifs, cette productivité supplémentaire pourrait faire augmenter les salaires des travailleurs utilisant l'IA à usage général qui subissent désormais une demande accrue pour leur travail. Deuxièmement, l'IA à usage général peut faire augmenter les salaires en stimulant la croissance économique et en stimulant la demande de main-d'œuvre dans des tâches qui ne sont pas encore automatisées (663, 664). Troisièmement, l'IA à usage général peut conduire à la création de tâches et de professions entièrement nouvelles pour les travailleurs (641, 644, 664). Cependant, l'IA à usage général peut également exercer une pression à la baisse sur les salaires des travailleurs de certaines professions. L'IA à usage général augmentant l'offre de certaines compétences sur le marché du travail, elle peut réduire la demande d'humains possédant ces mêmes compétences. Les travailleurs spécialisés dans des tâches pouvant être automatisées par l'IA à usage général peuvent donc être confrontés à une baisse de salaire ou à une perte d'emploi (643). Par exemple, une étude a révélé que quatre mois après la sortie de ChatGPT, celle-ci avait entraîné une baisse de 2 % du nombre d'offres d'emploi en ligne et une baisse de 5,2 % des revenus mensuels des rédacteurs sur la plateforme (645). L'impact sur les salaires dans un secteur donné dépend en grande partie de l'ampleur de la demande supplémentaire pour les services de ce secteur lorsque les coûts baissent en raison des gains de productivité générés par l'IA à usage général. En outre, la part des bénéfices générés par l'IA qui sont captés par

Les travailleurs dépendront de facteurs tels que les structures du marché et les politiques du travail dans les secteurs concernés, qui varient considérablement d'un pays à l'autre.

L'IA à usage général aura probablement l'impact le plus significatif à court terme sur les emplois qui consistent principalement en des tâches cognitives. Plusieurs études montrent que les capacités de l'IA à usage général se chevauchent avec les capacités nécessaires pour effectuer des tâches dans un large éventail d'emplois, les tâches cognitives étant les plus susceptibles d'être affectées (637*, 640, 665, 666). Les recherches ont également révélé que l'IA à usage général offre des gains de productivité importants aux travailleurs effectuant de nombreux types de tâches cognitives. Cela inclut le travail dans des professions telles que le conseil en stratégie (667), le travail juridique (668), la rédaction professionnelle (651), programmation informatique (113) et autres. Par exemple, les agents du service client ont bénéficié d'une augmentation moyenne de productivité de 14 % grâce à l'utilisation de l'IA à usage général (662). De plus, il a été constaté que les développeurs de logiciels exécutaient une tâche de codage illustratif 55,8 % plus rapidement lorsqu'ils avaient accès à un assistant de programmation d'IA à usage général (114*). Les secteurs qui dépendent fortement des tâches cognitives, tels que l'information, l'éducation et le secteur des services professionnels, scientifiques et techniques, adoptent également l'IA à des taux plus élevés, ce qui suggère que les travailleurs de ces secteurs sont sur le point d'être les plus touchés par l'IA à usage général à court terme (669).

Les agents d'IA ont le potentiel d'influencer les travailleurs de manière plus significative que les systèmes d'IA à usage général qui nécessitent une surveillance humaine importante. Les « agents d'IA » sont des systèmes d'IA à usage général qui peuvent accomplir des tâches en plusieurs étapes dans la poursuite d'un objectif de haut niveau avec peu ou pas de surveillance humaine. Cela signifie que les agents sont capables d'enchaîner plusieurs tâches complexes, automatisant potentiellement des flux de travail entiers plutôt que des tâches individuelles (670). En supprimant la nécessité d'une intervention humaine dans de longues séquences de travail, les agents d'IA pourraient effectuer des tâches et des projets à moindre coût que les agents d'IA.

Les systèmes d'IA à usage général nécessitent une plus grande surveillance humaine (671, 672). Cela devrait encourager l'adoption accrue d'agents à des fins d'automatisation dans des environnements économiquement compétitifs (671, 673). L'accélération de l'automatisation qui en résulterait pourrait entraîner une perturbation plus rapide des demandes de compétences et des salaires dans de nombreux secteurs (670), ce qui donnerait aux décideurs politiques moins de temps pour mettre en œuvre des mesures politiques qui renforcent la résilience des travailleurs.

La perte involontaire d'emploi peut causer des préjudices durables et graves aux travailleurs concernés. Des études montrent que les travailleurs licenciés subissent de fortes baisses de revenus et de consommation immédiatement après avoir été licenciés, les déficits de revenus persistant pendant des années par la suite (674, 675). Les estimations des baisses de salaire même après un réemploi varient de 5 à 30 % pendant 20 ans après le licenciement (676, 677, 678, 679). La perte involontaire d'emploi peut également affecter considérablement la santé physique, des données suggérant que le licenciement augmente le risque de mortalité de 50 à 100 % dans l'année suivant la séparation et de 10 à 15 % par an pendant les 20 années suivantes (680). Des études établissent également un lien entre la perte d'emploi et des taux plus élevés de dépression (681), de suicide (682), de maladies liées à l'alcool (682) et d'impacts négatifs sur le niveau d'éducation des enfants (683). Étant donné le potentiel de l'IA à usage général à provoquer des suppressions d'emplois, ces résultats soulignent l'importance des mesures politiques pour soutenir les travailleurs concernés.

L'amélioration des capacités de l'IA à usage général augmentera probablement les risques que les systèmes actuels présentent pour l'autonomie des travailleurs et le bien-être au travail. Les systèmes d'IA à usage général actuels sont déjà utilisés pour attribuer des tâches, surveiller la productivité et évaluer les performances des travailleurs dans des environnements allant des entrepôts aux centres d'appels (684). Bien que ces systèmes puissent augmenter la productivité (685), des études montrent qu'ils nuisent souvent au bien-être des travailleurs en raison d'une surveillance continue et de décisions relatives à la charge de travail basées sur l'IA (686). De nombreux employeurs adoptent ces systèmes sans effectuer suffisamment de tests ou sans comprendre pleinement leurs impacts sur la main-d'œuvre (687). Cela peut être particulièrement inquiétant lorsque les systèmes de gestion de l'IA influencent des décisions critiques telles que l'embauche et le licenciement (687). Il reste à voir si l'IA à usage général permettra une gestion algorithmique plus étendue que les systèmes d'IA à usage général qui sont souvent utilisés aujourd'hui. Si les systèmes d'IA à usage général améliorent l'intégration et l'analyse de divers flux de données, cela permettrait probablement une surveillance et une prise de décision plus granulaires sur les lieux de travail, augmentant potentiellement à la fois l'efficacité et les risques pour l'autonomie des travailleurs.

L'IA à usage général pourrait accroître les inégalités de revenus au sein des pays en offrant de plus grandes augmentations de productivité aux salariés à hauts revenus, mais les effets sont susceptibles de varier selon les pays. Au cours des dernières décennies, l'automatisation des tâches routinières a accru les inégalités salariales aux États-Unis en déplaçant les travailleurs à salaire moyen des emplois où ils bénéficiaient auparavant d'un avantage comparatif (688, 689, 690). Par exemple, une étude estime que 50 à 70 % de l'augmentation des inégalités salariales aux États-Unis au cours des quatre dernières décennies peut s'expliquer par la baisse relative des salaires des travailleurs spécialisés dans les tâches routinières dans les secteurs qui ont connu des niveaux élevés d'automatisation (688).

L'IA à usage général pourrait concurrencer les travailleurs humains de la même manière, en réduisant potentiellement les salaires de certains travailleurs (691, 692) tout en étant plus susceptible d'augmenter la productivité de ceux qui occupent déjà des emplois relativement bien rémunérés (voir la figure 2.6) (637*). Une simulation suggère que l'IA pourrait accroître de 10 % les inégalités salariales entre les professions à revenu élevé et à faible revenu dans un certain nombre de pays.

décennie dans les économies avancées (640). Cependant, pour de nombreux types de tâches cognitives, il existe des preuves qu'au niveau actuel des capacités du modèle, ceux qui ont moins d'expérience ou des compétences plus élémentaires obtiennent souvent les plus grandes augmentations de productivité en utilisant l'IA à usage général (114*, 651, 662, 667, 668). Cela suggère que dans les professions axées sur les tâches cognitives, les travailleurs moins bien payés pourraient en fait obtenir une plus grande augmentation que les hauts revenus et que les inégalités salariales au sein de ces professions pourraient diminuer (693). La manière dont ces effets compensatoires se manifesteront dans l'ensemble de l'économie est incertaine et est susceptible de varier selon les pays, les secteurs et les professions.

Exposition aux grands modèles linguistiques (LLM) par revenu

Part de toutes les tâches au sein des professions qui sont exposées aux LLM et aux LLM partiels

Les logiciels sont présentés par rapport au salaire annuel médian de la profession. Les données reflètent l'évaluation humaine.

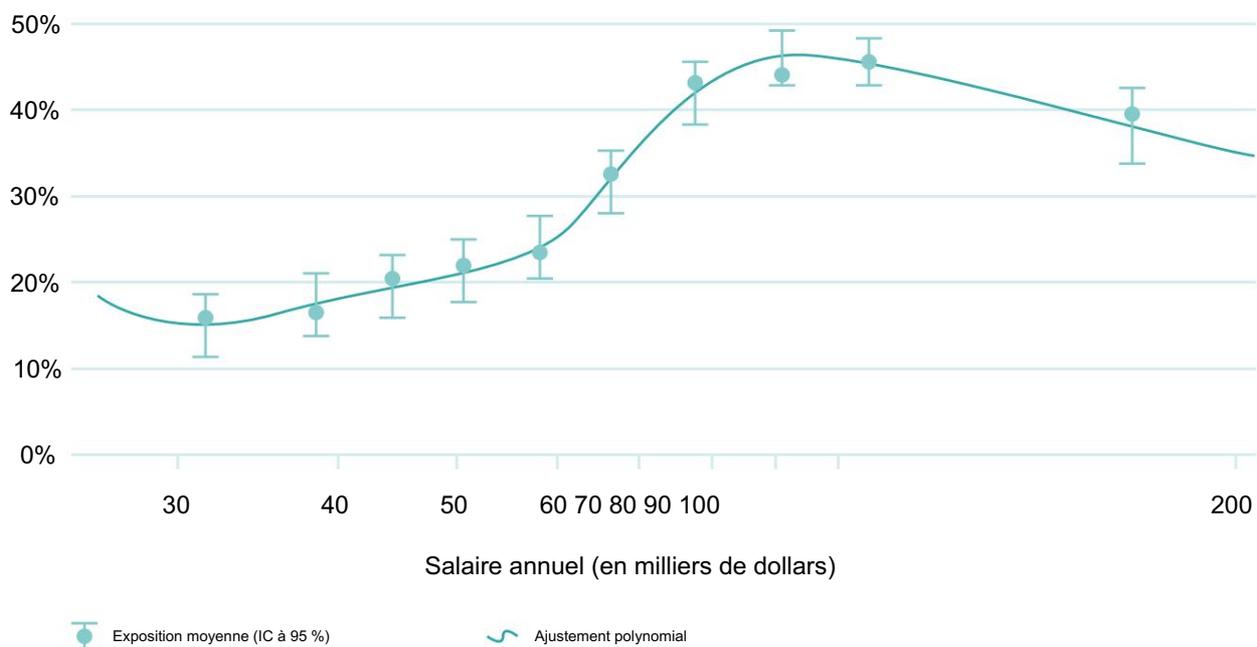


Figure 2.6 : Les grands modèles linguistiques (LLM) ont un impact économique inégal sur les différentes parties de la distribution des revenus. L'exposition est la plus élevée pour les tâches des travailleurs se situant à l'extrémité supérieure des salaires annuels, culminant à environ 90 000 \$/an aux États-Unis, tandis que les revenus faibles et moyens sont nettement moins exposés. Dans cette figure, « exposition » signifie le potentiel de gains de productivité de l'IA, qui peut se manifester par une augmentation du nombre de travailleurs et des augmentations de salaires ou par une automatisation et des baisses de salaires, en fonction de divers autres facteurs. Source : Eloundou et al., 2024 (637*).

L'automatisation du travail par l'IA à des fins générales risque d'exacerber les inégalités en réduisant la part des revenus qui revient aux travailleurs par rapport aux propriétaires de capital. À l'échelle mondiale, la part des revenus du travail a diminué d'environ six points de pourcentage entre 1980 et 2022 (694). En règle générale, 10 % de tous les salariés reçoivent la majorité des revenus du capital (695, 696). Si l'IA automatise une part importante du travail, ces tendances pourraient s'intensifier en réduisant les opportunités de travail pour les salariés et en augmentant le rendement de la propriété du capital (697, 698). En outre, des données suggèrent que l'IA à des fins générales peut aider à la création de grandes entreprises « superstars » qui captent une part importante des bénéfices économiques, ce qui accroîtrait encore davantage l'inégalité entre le capital et le travail (699).

Les technologies d'IA à usage général risquent d'aggraver les inégalités mondiales si elles sont principalement adoptées par les pays à revenu élevé (PRE). Les PRE comptent une part plus élevée d'emplois axés sur les tâches cognitives qui sont les plus exposés aux impacts de l'IA à usage général (640). Ces pays disposent d'infrastructures numériques plus solides, d'une main-d'œuvre qualifiée et d'écosystèmes d'innovation plus développés (700) (voir 2.3.2.

[Français Fracture mondiale de la R&D en IA](#)). Cela les place en position de capter les gains de productivité de l'IA à usage général plus rapidement que les marchés émergents et les économies en développement. Cela contribuerait à des trajectoires de croissance des revenus divergentes et à un écart croissant entre les pays à revenu élevé et les pays à revenu faible et intermédiaire (PRFI) (701). Si l'IA la plus avancée et la plus automatisée du travail est utilisée par les entreprises des pays à revenu élevé, cela pourrait également attirer des investissements en capital supplémentaires dans ces pays et accentuer encore davantage la divergence économique entre les régions à revenu élevé et à faible revenu (702). En outre, à mesure que les entreprises des économies avancées adoptent l'IA à usage général, elles pourraient trouver plus rentable d'automatiser la production au niveau national plutôt que de délocaliser le travail, érodant ainsi une voie de développement traditionnelle pour les économies en développement qui exportent des services à forte intensité de main-d'œuvre (703). Une étude suggère que cette dynamique est plus susceptible de se jouer dans les pays où une grande partie de la main-d'œuvre travaille dans des services informatiques externalisés tels que le service client, la rédaction et les emplois de l'économie numérique à la demande (704). Cependant, l'impact précis sur les marchés du travail dans les économies en développement reste flou. D'une part, ces économies pourraient être confrontées à un double défi : perdre des emplois existants au profit de l'automatisation tout en ayant plus de mal à attirer de nouveaux investissements, les avantages liés au coût de la main-d'œuvre devenant moins pertinents. D'autre part, si l'IA à usage général est largement adoptée dans les économies en développement, elle pourrait permettre d'accroître la productivité de certains travailleurs qualifiés (662, 705, 706), ce qui pourrait leur permettre de rivaliser pour des opportunités de travail à distance avec leurs homologues mieux rémunérés dans les pays à revenu élevé.

Français Depuis la publication du rapport intermédiaire, de nouvelles données suggèrent que les taux d'adoption de l'IA à usage général par les individus pourraient être plus rapides que les technologies précédentes telles qu'Internet ou les ordinateurs personnels, bien que le rythme d'adoption par les entreprises varie considérablement selon les secteurs (voir la figure 2.7) (661). Par exemple, une enquête récente menée aux États-Unis a révélé que plus de 24 % des travailleurs utilisent l'IA générative au moins une fois par semaine, et un sur neuf l'utilise quotidiennement au travail (661). Les taux d'adoption des entreprises varient considérablement selon les secteurs (707). Par exemple, aux États-Unis, environ 18,1 % des entreprises du secteur de l'information déclarent utiliser l'IA (au sens large), tandis que seulement 1,4 % dans la construction et l'agriculture le font (669). Parmi les entreprises qui déclarent utiliser l'IA, 27 % déclarent remplacer les tâches des travailleurs, tandis que seulement 5 % signalent des changements d'emploi dus à l'IA, dont plus de la moitié sont des augmentations d'emploi plutôt que des diminutions (708). Les données actuelles sur les taux d'adoption de l'IA à usage général sont limitées par la collecte limitée de données internationales, en particulier en dehors des États-Unis, bien qu'une enquête menée auprès de plus de 15 000 travailleurs dans 16 pays ait révélé que 55 % des répondants utilisent l'IA générative au moins une fois par semaine dans leur travail (709).

Partout dans le monde, il existe un écart important entre les sexes en termes d'adoption et d'impact potentiel de l'IA à usage général sur le marché du travail. Par exemple, une méta-analyse récente de dix études menées dans différents pays suggère que les femmes sont 24,6 % moins susceptibles d'utiliser l'IA générative que les hommes (710).

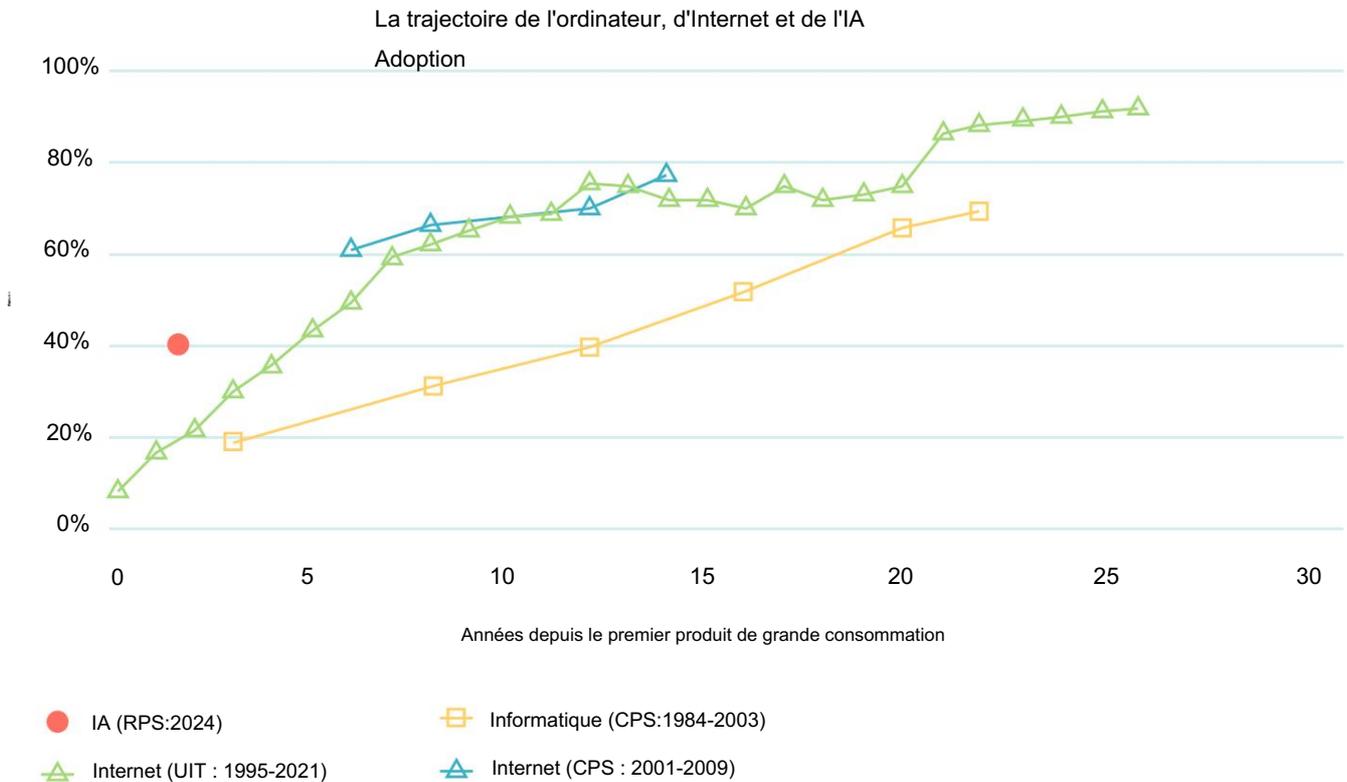


Figure 2.7 : Jusqu'à présent, l'IA générative semble avoir été adoptée à un rythme plus rapide que les PC ou Internet aux États-Unis. Cette adoption plus rapide par rapport aux PC est due à une utilisation beaucoup plus importante en dehors du travail, probablement en raison de différences de portabilité et de coût. Source : Bick et al., 2024 (661).

De plus, depuis la publication du rapport intermédiaire, de nouvelles preuves ont montré que l'IA à usage général peut générer des gains de productivité significatifs dans des environnements de travail réels et est sur le point de générer des gains dans les domaines de la science et de la R&D. De nouvelles preuves démontrent les impacts sur la productivité dans plusieurs environnements de travail réels et constatent que ces impacts varient selon la profession, selon l'entreprise et sont influencés par le taux d'adoption et d'utilisation au sein d'une entreprise (711*). Des recherches récentes ont également révélé que chaque multiplication par 10 de la puissance de calcul utilisée pour former un modèle permettait aux travailleurs d'effectuer certaines tâches de traduction 12,3 % plus rapidement et avec une qualité améliorée lorsqu'ils utilisaient le modèle comme assistant (706). La mesure dans laquelle cette relation s'applique à d'autres tâches professionnelles reste toutefois incertaine. En outre, l'impact des gains de productivité de l'IA à usage général sur le développement des compétences des travailleurs par rapport au déclin de leurs compétences est un sujet de recherche émergent. Une étude récente a révélé que si l'utilisation de ChatGPT peut améliorer certaines capacités des travailleurs, les compétences et les connaissances ne sont généralement pas conservées une fois l'accès supprimé (712*).

Enfin, une étude récente sur le marché du travail américain a montré que les emplois axés sur la technologie et la R&D étaient ceux qui affichaient la plus forte proportion de tâches exposées à des gains de productivité potentiels grâce à l'IA à usage général (637*). Cela suggère que les impacts importants sur la productivité récemment observés des systèmes d'IA à usage restreint sur la découverte scientifique (713) pourraient potentiellement être accentués par l'IA à usage général dans un éventail plus large de tâches de R&D. Cela est remarquable, car l'augmentation de la productivité de la R&D peut considérablement stimuler le progrès technologique et la croissance économique (714).

Les principales lacunes en matière de données probantes concernant les risques sur le marché du travail comprennent des impacts incertains sur l'emploi à long terme, des données internationales limitées sur l'adoption et des réponses politiques non testées. Il n'existe pas d'études exhaustives sur les effets à long terme de l'IA à usage général sur l'emploi et les salaires dans différents secteurs. Les modèles d'adoption sont mal compris en dehors des États-Unis, ce qui rend difficile d'anticiper les impacts internationaux. Les données sur la création de nouveaux emplois grâce à l'adoption de l'IA à usage général sont insuffisantes pour guider les programmes de reconversion des travailleurs. Plus important encore pour les décideurs politiques, il existe peu de preuves sur les interventions qui protègent efficacement les travailleurs pendant les transitions technologiques. Par exemple, bien que la reconversion soit souvent suggérée comme une réponse à l'évolution des demandes de compétences, il existe peu de preuves de son efficacité (715), en particulier compte tenu de la rapidité avec laquelle l'IA à usage général pourrait modifier les compétences requises sur le lieu de travail. Ces lacunes existent en partie parce que l'IA à usage général est encore naissante, ce qui rend les impacts à long terme difficiles à mesurer, et en partie parce qu'il est difficile d'isoler les effets de l'IA à usage général des autres facteurs économiques. Le rythme rapide du développement de l'IA à usage général signifie également que les preuves de l'efficacité des réponses politiques aux changements technologiques antérieurs peuvent ne pas être transposées dans ce contexte.

Pour les décideurs politiques qui travaillent sur les risques liés à l'IA à usage général sur le marché du travail, les principaux défis consistent à trouver un équilibre entre l'innovation en matière d'IA et la protection des travailleurs, à créer des politiques qui peuvent s'adapter rapidement à l'évolution des impacts et à garantir que les avantages économiques sont partagés à l'intérieur des pays et entre eux. L'un des principaux défis consiste à trouver un équilibre entre l'innovation (qui pourrait stimuler la productivité et la croissance) et la protection des travailleurs contre les baisses de salaires et les préjudices associés à la perte involontaire d'emploi (654). Les politiques doivent être adaptables compte tenu de la rapidité du développement de l'IA à usage général et de l'incertitude quant aux impacts futurs (716). Cependant, les décideurs politiques sont confrontés au défi de définir des politiques flexibles tout en offrant suffisamment de certitude réglementaire pour faciliter les décisions d'investissement des entreprises et de formation des travailleurs. Un suivi étroit des principales tendances peut aider les décideurs politiques à mieux anticiper les impacts de l'IA à usage général sur le marché du travail. Il s'agit notamment des taux d'adoption de l'IA par secteur, des changements dans la répartition des salaires entre les industries, de l'émergence de nouvelles catégories d'emplois et des changements dans les demandes de compétences des employeurs à mesure que les systèmes d'IA à usage général progressent et sont plus largement adoptés. Enfin, étant donné que les impacts de l'IA sur le marché du travail devraient varier considérablement d'un pays à l'autre, les décideurs politiques sont confrontés à des défis pour coordonner les réponses internationales afin d'éviter un élargissement de la fracture économique mondiale et de garantir que l'IA puisse accélérer la croissance économique inclusive.

2.3.2. Fracture mondiale de la recherche et du développement en IA

INFORMATIONS CLÉS

- Les grandes entreprises des pays dotés d'infrastructures numériques solides sont les chefs de file de la recherche et du développement en matière d'IA à usage général, ce qui pourrait entraîner une augmentation des inégalités et des dépendances à l'échelle mondiale. Par exemple, en 2023, la majorité des modèles d'IA à usage général notables (56 %) ont été développés aux États-Unis. Cette disparité expose de nombreux pays à revenu faible ou intermédiaire à des risques de dépendance et pourrait exacerber les inégalités existantes.
- L'augmentation des coûts de développement de l'IA à usage général est la principale raison de cette « fracture de la R&D en IA ». L'accès à des quantités importantes et coûteuses de puissance de calcul est devenu une condition préalable au développement d'une IA polyvalente avancée. Les institutions universitaires et la plupart des entreprises, en particulier celles des pays à revenu faible ou intermédiaire, n'ont pas les moyens de rivaliser avec les grandes entreprises technologiques.
- Les tentatives visant à combler le fossé en matière de recherche et développement en IA n'ont pas été couronnées de succès. Un nombre croissant d'efforts ont été consacrés à la démocratisation de l'accès à l'informatique, à l'investissement dans la formation aux compétences en IA dans les pays à revenu faible et intermédiaire et à l'open source des principaux modèles d'IA. Mais ces efforts nécessiteront des investissements financiers considérables et un temps de mise en œuvre important.
- Des travaux récents suggèrent que le fossé en matière de R&D en IA pourrait encore se creuser en raison d'une tendance à l'augmentation des coûts de R&D à la frontière. Depuis la publication du rapport intermédiaire (mai 2024), les chercheurs ont publié de nouvelles preuves sur l'augmentation des coûts de développement d'une IA de pointe, les disparités croissantes dans la concentration des talents en IA et la centralisation croissante des ressources informatiques nécessaires à la formation de grands modèles d'IA à usage général. • Il existe un manque de preuves sur l'efficacité des moyens potentiels de répondre à la R&D en IA
Par exemple, l'impact des programmes de formation à l'IA ou des investissements dans les infrastructures dans les pays à revenu faible ou intermédiaire reste flou.

Définitions clés

- Fracture numérique : disparité dans l'accès aux technologies de l'information et de la communication (TIC), notamment à Internet, entre différentes régions géographiques ou différents groupes de personnes. • Fracture en R&D en IA : disparité dans la recherche et le développement en IA entre différentes régions géographiques, causée par divers facteurs, notamment une répartition inégale de la puissance de calcul, des talents, des ressources financières et des infrastructures.
- Infrastructure numérique : les services et installations fondamentaux nécessaires à la numérisation technologies de fonctionnement, y compris le matériel, les logiciels, les réseaux, les centres de données et les systèmes de communication.
- Travail fantôme : travail caché effectué par les travailleurs pour soutenir le développement et le déploiement de modèles ou de systèmes d'IA (par exemple via l'étiquetage des données).

Français La répartition mondiale inégale des ressources informatiques, des talents, des ressources financières et de l'infrastructure numérique contribue à une fracture en matière de R&D en IA qui pourrait exposer de nombreux pays à revenu faible et intermédiaire à des risques de dépendance et entraver leur progression dans la R&D en IA à usage général. Les coûts financiers élevés du développement et de l'exploitation de systèmes d'IA à usage général (27) peuvent limiter la production de R&D en IA à usage général des pays à revenu faible et intermédiaire, ce qui pourrait exacerber les inégalités existantes. Les chercheurs des pays à revenu faible et intermédiaire, qui ne sont souvent pas en mesure de former des LLM en raison des coûts élevés, auront tendance à s'appuyer sur les modèles à pondération ouverte existants, qui sont principalement développés dans les pays dotés d'une infrastructure numérique solide (717). Ces modèles ne sont probablement pas en mesure de saisir pleinement les nuances (structure grammaticale, écritures non latines, différences tonales, etc.) des langues non occidentales, qui sont sous-représentées dans les données de formation, ce qui entraîne une précision moindre (718).

De plus, la dépendance vis-à-vis des entreprises nord-américaines et chinoises pour l'accès aux calculs et aux modèles ouverts s'accompagne généralement de restrictions en matière de droits d'auteur et de confidentialité qui limitent la capacité des chercheurs et des développeurs de nombreux pays à revenu faible ou intermédiaire à créer des modèles de pointe (719). Ainsi, ces chercheurs dépendent souvent de collaborations avec des parties prenantes de pays dotés d'infrastructures numériques plus solides pour accéder aux calculs et publier dans des lieux de premier plan. Enfin, alors que des pays comme les États-Unis et la Chine continuent d'être les chefs de file de la production de talents qualifiés en IA, les chercheurs et les étudiants d'autres pays dépendent souvent des institutions de ces pays leaders pour leur avancement académique et professionnel en IA. Cela peut exacerber les disparités dans la R&D en IA, car les talents se déplacent d'autres pays vers des pays où une industrie de l'IA est déjà concentrée (720).

L'un des principaux facteurs de fracture en matière de R&D en IA est la différence d'accès aux ressources informatiques entre les différents acteurs. Cela inclut l'accès inégal aux ressources informatiques puissantes (unités de traitement graphique (GPU), centres de données, services cloud, etc.) qui sont nécessaires pour former et déployer des modèles d'IA de grande taille et complexes. Ces dernières années, ce fossé s'est creusé (721, 722). Cette inégalité d'accès est particulièrement apparente dans la mesure dans laquelle les grandes entreprises d'IA et les laboratoires universitaires d'IA ont accès aux ressources informatiques. Les estimations montrent que les entreprises technologiques américaines sont les principaux acheteurs de GPU NVIDIA H100, l'un des types de puces GPU les plus puissants du marché, explicitement conçu pour l'IA (723).

Cependant, plusieurs grandes entreprises technologiques ont toutes récemment annoncé qu'elles développaient des puces d'IA personnalisées pour réduire leur dépendance à la chaîne d'approvisionnement en puces d'IA, ouvrant ainsi la voie à un accès plus large aux GPU. Mais le coût exceptionnellement élevé des GPU (généralement 20 000 à 30 000 dollars pour les GPU haut de gamme tels que le H100 en novembre 2024), dont des milliers ou des dizaines de milliers sont généralement utilisés pour former un modèle d'IA polyvalent de premier plan, pourrait encore empêcher la plupart des pays à revenu faible ou intermédiaire de se permettre ce niveau d'infrastructure d'IA. Le coût croissant de la création et de la maintenance des centres de données contribue également à l'inégalité d'accès au calcul. Au cours de la dernière décennie, les grandes entreprises technologiques ont augmenté leurs investissements dans les centres de données, Google dévoilant un centre de données de 600 millions de dollars dans le Nebraska en 2022 (724*) et annonçant récemment un plan de construction d'un centre de données d'un milliard de dollars dans le Missouri (725). Meta a investi plus de 2 milliards de dollars dans un centre de données en Oregon (726*), et Microsoft a annoncé une initiative d'un milliard de dollars pour construire un campus de centres de données, ainsi que d'autres efforts de développement de l'IA, au Kenya (727*). Si de tels efforts contribueront considérablement à accroître l'accès à l'informatique en général, il est peu probable qu'ils réduisent de manière significative la fracture en matière de recherche et développement en IA.

Les disparités dans la concentration des talents qualifiés contribuent également à la fracture mondiale en matière de R&D en IA. La R&D en IA est principalement concentrée dans deux pays – les États-Unis et la Chine – qui ont réalisé des investissements importants pour recruter et retenir les talents en IA. Les États-Unis comptent le plus grand pourcentage de chercheurs d'élite en IA, abritent la majorité des institutions qui mènent des recherches de haut niveau et sont la première destination des talents en IA à l'échelle mondiale (728). En outre, il existe des disparités dans l'accès des étudiants aux programmes de diplômes liés à l'IA, car de nombreuses universités de premier plan en IA sont basées aux États-Unis ou au Royaume-Uni (729), et la grande majorité des cours universitaires en anglais sur l'IA sont proposés au Royaume-Uni, aux États-Unis et au Canada (730). Alors que certains PRFI, comme l'Inde et la Malaisie, augmentent leur offre de cours en IA (731), il reste encore beaucoup à faire pour comprendre cette disparité, car les données sur les programmes universitaires formels en IA dans les PRFI sont limitées, en particulier ceux proposés dans des langues autres que l'anglais.

La délégation de tâches d'IA de niveau inférieur aux travailleurs des pays à revenu faible et intermédiaire a donné naissance à une industrie du « travail fantôme ». La demande croissante de données pour former des systèmes d'IA à usage général, y compris des retours humains pour aider à la formation, a encore accru le recours au « travail fantôme » (732). Le « travail fantôme » est principalement un travail caché effectué par des travailleurs – souvent dans des conditions précaires – pour soutenir le développement de modèles d'IA. Des entreprises ont vu le jour pour aider les grandes entreprises technologiques à externaliser divers aspects de la production de données, notamment la collecte, le nettoyage et l'annotation des données. Ce travail peut offrir des opportunités aux personnes des pays à revenu faible et intermédiaire. D'un autre côté, la nature contractuelle de ce travail offre souvent peu d'avantages et de protections aux travailleurs et moins de stabilité de l'emploi, car les plateformes font tourner les marchés pour trouver une main-d'œuvre moins chère. Des recherches ont montré que ces travailleurs sont exposés à des contenus explicites, à des horaires irréguliers, à de lourdes charges de travail et à une mobilité sociale et économique limitée (733, 734, 735, 736).

L'exposition à un contenu aussi graphique peut entraîner un syndrome de stress post-traumatique et d'autres traumatismes mentaux (737, 738).

Depuis la publication du rapport intermédiaire, de plus en plus de preuves de l'augmentation des coûts associés au développement de l'IA à usage général sont apparues, ce qui laisse penser que le fossé en matière de R&D en IA se creusera encore davantage. Le développement de modèles d'IA à usage général notables est toujours mené par des entreprises de pays dotés d'une solide infrastructure numérique et d'un accès au calcul, et les capacités de ces modèles augmentent. Les chercheurs ont fourni des preuves solides que l'utilisation de ressources telles que l'électricité dans le développement de l'IA est en augmentation (739). Le coût de la formation de modèles d'IA de pointe a augmenté de 2 à 3 fois par an au cours des huit dernières années et pourrait atteindre un coût de plus d'un milliard de dollars d'ici 2027 (27). Cependant, certains éléments indiquent une amélioration de la concentration des talents et du développement de modèles de pointe dans les pays à revenu faible et intermédiaire.

L'Inde, par exemple, a particulièrement bien réussi à accroître sa concentration de talents qualifiés en IA, qui a augmenté de 263 % depuis 2016 (740). Les recherches indiquent que le développement de l'IA à usage général peut avoir un impact significatif sur les services informatiques externalisés vers les pays à revenu faible ou intermédiaire, tels que le service client, la rédaction et le travail à la demande (704).

L'une des principales lacunes en matière de recherche et développement en IA est le manque de preuves sur les solutions réalisables.

Les grandes entreprises technologiques investissent de plus en plus dans l'IA et dans la formation aux compétences numériques en Afrique, en Amérique latine et en Asie, et ces programmes sont susceptibles d'augmenter à mesure que ces régions élargissent les capacités des modèles de pointe pour les consommateurs locaux.

Cependant, rien ne prouve que cette formation améliore la production de modèles d'IA significatifs, en particulier dans les pays à revenu faible et intermédiaire. Il n'existe également que peu de preuves des avantages des investissements dans les infrastructures spécifiques à l'IA, compte tenu des grandes disparités de talents en IA entre de nombreux pays à revenu faible et intermédiaire et des pays comme les États-Unis et la Chine. À l'heure actuelle, il n'est pas certain que l'accès à l'infrastructure augmenterait les talents, ou que cette infrastructure resterait inutilisée en raison d'un manque d'experts qualifiés. Il existe également peu de données sur l'ampleur de la fracture de la R&D en IA, car les indicateurs mesurent souvent les résultats de la recherche dans des revues et des conférences de premier plan, qui sont toutes publiées en anglais. Des obstacles structurels, tels que les restrictions de visa et les charges financières, empêchent souvent les chercheurs internationaux qualifiés, en particulier dans les pays à revenu faible et intermédiaire, d'assister à des conférences majeures ou de publier dans des revues coûteuses. Les effets de la fracture de la R&D en IA sont également des retombées de la fracture numérique existante (741), ce qui rend difficile de distinguer les impacts spécifiques de l'IA à usage général sur la fracture mondiale de la R&D en IA.

La réduction de la fracture en matière de recherche et développement en IA est un problème difficile à résoudre pour les décideurs politiques. Les coûts de développement de l'IA à usage général sont inaccessibles pour la majorité des PRFI, et les investissements dans les infrastructures de base telles que les réseaux électriques et les réseaux Internet sont estimés à des milliards de dollars (USD) pour des pays comme le Nigéria (742). En outre, aucun de ces pays ne dispose d'entreprises capables de supporter individuellement les dépenses liées au développement de systèmes d'IA à usage général. Les données probantes sur les résultats de la formation aux compétences numériques sont limitées, ce qui peut entraver les efforts supplémentaires visant à développer des programmes de formation aux compétences ciblées susceptibles d'avoir un impact significatif sur les contributions des PRFI aux modèles d'IA à usage général. Il existe également des projections selon lesquelles les disparités dans la concentration des talents en IA pourraient s'accroître. Des pays comme les États-Unis, le Royaume-Uni, la Chine et des pays d'Europe augmentent rapidement le recrutement de talents en IA, certains offrant des voies d'immigration aux talents qualifiés pour contribuer à la R&D en IA dans leurs pays respectifs (743, 744, 745). Les décideurs politiques, en particulier dans de nombreux PRFI, devront analyser les implications de cette situation pour leur autonomie régionale et leurs efforts pour réduire la fracture en matière de recherche et développement en IA.

2.3.3. Concentration du marché et points de défaillance uniques

INFORMATIONS CLÉS

- Les parts de marché de l'IA à usage général ont tendance à être très concentrées entre quelques acteurs, ce qui peut créer une vulnérabilité aux défaillances systémiques. Le degré élevé de concentration du marché peut conférer à un petit nombre de grandes entreprises technologiques un pouvoir considérable sur le développement et le déploiement de l'IA, ce qui soulève des questions sur leur gouvernance. L'utilisation généralisée de quelques modèles d'IA à usage général peut également rendre les secteurs de la finance, de la santé et d'autres secteurs critiques vulnérables aux défaillances systémiques si l'un de ces modèles présente des problèmes.
- Le marché est tellement concentré en raison des barrières élevées à l'entrée. Les modèles d'IA polyvalents de pointe nécessitent un investissement initial considérable. Par exemple, le coût global du développement d'un modèle de pointe peut actuellement atteindre des centaines de millions de dollars américains. Les principaux facteurs de coût sont la puissance de calcul, la main-d'œuvre hautement qualifiée et les vastes ensembles de données.
- De plus, les leaders du marché bénéficient d'une dynamique auto-renforçante qui récompense les gagnants. Les économies d'échelle permettent aux plus grandes entreprises d'IA de répartir les coûts de développement ponctuels sur une base de clients toujours plus large, créant ainsi un avantage de coût par rapport aux plus petites entreprises. Les effets de réseau permettent en outre aux plus grandes entreprises de former de futurs modèles avec des données utilisateur générées par des modèles plus anciens.
- La concentration du marché a continué de persister en 2024. Depuis la publication du rapport intermédiaire (mai 2024), le consensus précédent selon lequel la concentration du marché de l'IA à usage général est élevée continue de se maintenir.
- Il existe peu de recherches sur la prévision ou l'atténuation des points de défaillance uniques de l'IA. Cela crée des difficultés pour les décideurs politiques. L'absence de méthodes de prévision fiables sur la manière dont les défaillances peuvent se propager dans les systèmes interconnectés rend ces risques difficiles à évaluer.

Définitions clés

- Concentration du marché : degré auquel un petit nombre d'entreprises contrôlent un marché
 - Point de défaillance unique : un élément d'un système plus vaste dont la défaillance perturbe l'ensemble du système. Par exemple, si un seul système d'IA joue un rôle central dans l'économie ou dans une infrastructure critique, son dysfonctionnement pourrait provoquer des perturbations généralisées dans toute la société.

Le développement d'une IA polyvalente de pointe nécessite d'énormes investissements financiers, atteignant souvent des centaines de millions de dollars américains (voir la figure 2.8). Ces coûts proviennent principalement de trois domaines clés : des ressources informatiques spécialisées, une expertise hautement qualifiée en IA et l'accès à de vastes ensembles de données, qui sont souvent propriétaires et coûteux. Les ressources informatiques comprennent du matériel avancé tel que les GPU (unités de traitement graphique) et les TPU (unités de traitement tensoriel),

L'infrastructure cloud et l'énergie nécessaire à la formation de modèles d'IA à grande échelle (739) sont également des facteurs importants de la création d'ensembles de données de haute qualité, en raison de processus tels que la collecte, l'annotation et le nettoyage (746, 747). En outre, le recrutement et la rétention de chercheurs, d'ingénieurs et de data scientists de haut niveau en IA sont très compétitifs et coûteux, car leur expertise est essentielle au développement d'algorithmes et d'architectures de pointe.

Les coûts estimés de formation des modèles d'IA ont fortement augmenté récemment

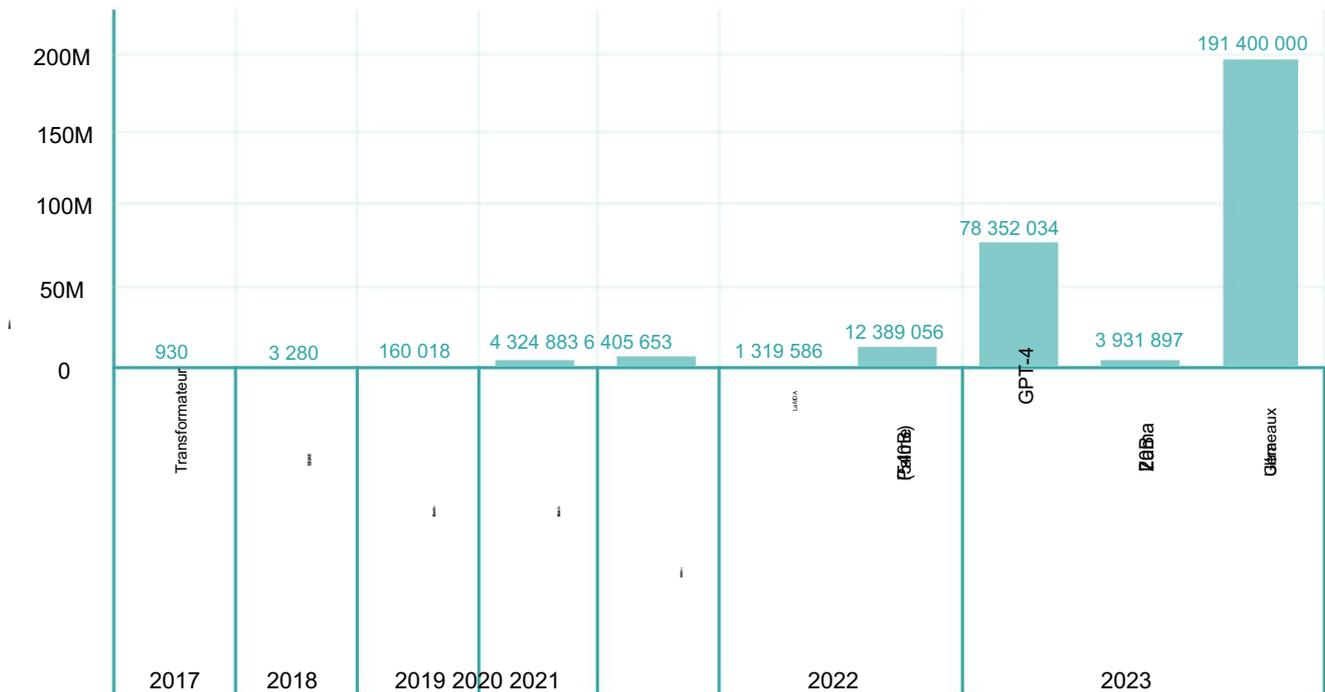


Figure 2.8 : Les coûts estimés de formation des modèles d'IA ont fortement augmenté au cours des dernières années. Seules quelques entreprises peuvent se permettre de former des modèles à un coût aussi élevé, ce qui accroît encore la concentration du marché. Source : Maslej et al., 2024a (730).

L'accès à des ensembles de données massifs est essentiel pour former des modèles d'IA de haute performance. Ces ensembles de données sont souvent propriétaires, ce qui confère aux entreprises établies un avantage concurrentiel (voir 1.3. [Capacités dans les années à venir](#)). Les grandes entreprises technologiques sont particulièrement bien placées pour surmonter ces obstacles en raison de leurs ressources financières, de leur infrastructure et de la propriété ou du contrôle qu'elles détiennent sur de vastes volumes de données via leurs plateformes et services existants. En revanche, les nouvelles entreprises sont confrontées à des obstacles importants pour acquérir les ensembles de données et la puissance de calcul nécessaires, ce qui entraîne une barrière à l'entrée élevée (74, 748, 749, 750). En conséquence, les petites entreprises sont souvent incapables de rivaliser, ce qui renforce la concentration du pouvoir de marché entre quelques acteurs dominants dans le secteur de l'IA.

Les systèmes d'IA à usage général bénéficient considérablement des économies d'échelle, car les modèles plus grands et plus gourmands en calcul ont tendance à surpasser leurs homologues plus petits sur de nombreux paramètres.

Les modèles à grande échelle, tels que ceux utilisés pour le traitement du langage naturel, la reconnaissance d'images et la prise de décision, sont capables de gérer un plus large éventail de tâches en raison de leur capacité accrue à traiter et à analyser de vastes quantités de données. À mesure que ces modèles prennent de l'ampleur, cela peut également entraîner

Une meilleure généralisation et une meilleure précision (751) renforcent la demande de systèmes d'IA polyvalents et hautement performants dans tous les secteurs. Cela crée une boucle de rétroaction dans laquelle les modèles à grande échelle, qui nécessitent des ressources informatiques importantes pour se développer, deviennent plus précieux et recherchés en raison de leurs performances et de leur polyvalence. Étant donné que les systèmes d'IA nécessitent des investissements initiaux importants en termes d'infrastructure et de développement (27), mais seulement de faibles coûts par requête, le coût moyen par utilisateur diminue à mesure que le système d'IA est fourni à davantage d'utilisateurs, ce qui reflète les économies d'échelle. Cela confère aux grandes entreprises un avantage concurrentiel, car elles peuvent répartir les coûts de développement sur une base de clients plus large, ce qui rend la concurrence difficile pour les petites entreprises. De plus, ces systèmes bénéficient d'effets de réseau : à mesure que davantage d'utilisateurs interagissent avec eux, ils génèrent de grandes quantités de nouvelles données qui peuvent être utilisées pour recycler et affiner les modèles (752, 753). Cet afflux constant de données générées par les utilisateurs améliore les performances des modèles, les rendant encore plus précieux et efficaces au fil du temps.

Ces tendances à la concentration du marché signifient que quelques entreprises domineront probablement la prise de décision concernant le développement et le déploiement de l'IA à usage général. Étant donné que la société dans son ensemble pourrait bénéficier autant que souffrir des décisions de ces entreprises, cela soulève des questions sur la gouvernance appropriée de ces quelques systèmes à grande échelle. Un seul modèle d'IA à usage général pourrait potentiellement influencer la prise de décision dans de nombreuses organisations et secteurs (748) de manières qui pourraient être bénignes, subtiles, involontaires ou délibérément exploitées. Il existe un risque d'utilisation malveillante de l'IA à usage général comme un puissant outil de manipulation, de persuasion, de censure et de contrôle par quelques entreprises ou gouvernements.

Part de marché mondiale des principaux fournisseurs de services d'infrastructure cloud au premier trimestre 2024

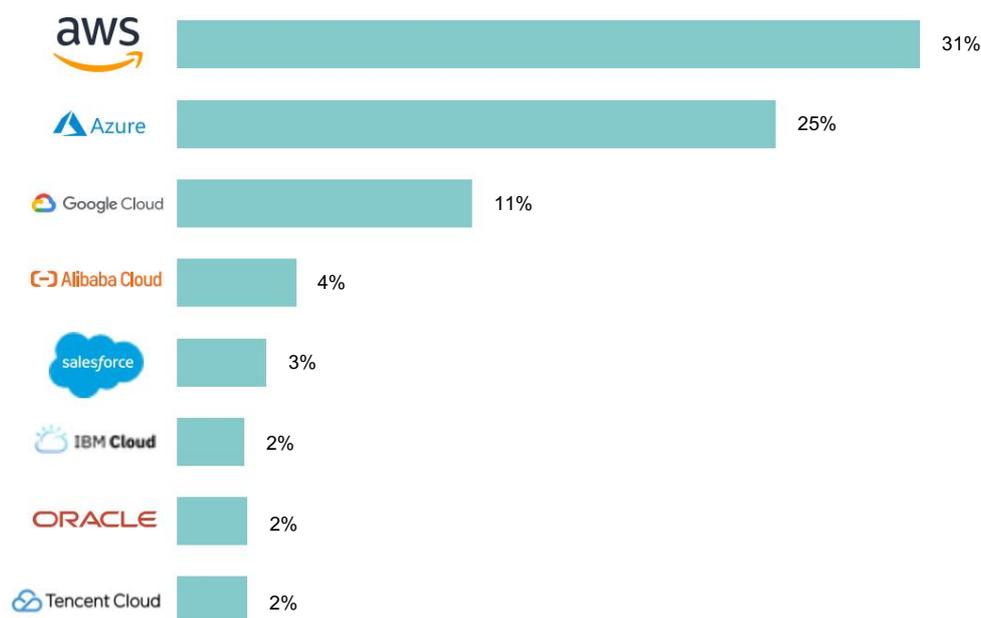


Figure 2.9 : Amazon (AWS), Microsoft (Azure) et Google contrôlent ensemble la concentration $\frac{3}{5}$ des services mondiaux de cloud computing, du pouvoir sur l'infrastructure essentielle de formation et de déploiement de l'IA dans seulement trois entreprises. Source : Richter, 2024 (756).

Français Depuis la publication du rapport intermédiaire, le consensus précédent selon lequel la concentration du marché de l'IA à usage général est élevée continue de se maintenir, et certaines nouvelles recherches suggèrent une dépendance croissante à l'égard des grandes entreprises d'IA. La dépendance croissante à l'égard des grandes entreprises technologiques pour l'accès au matériel essentiel (GPU), aux interfaces de modèles d'IA (API) et aux services de stockage en nuage a des implications importantes pour l'écosystème de l'IA (754). Trois entreprises seulement contrôlent 67 % des services de cloud computing (voir la figure 2.9). Cette dépendance consolide le pouvoir au sein de quelques acteurs majeurs, limitant la concurrence et l'innovation entre les petites entreprises qui n'ont pas les ressources pour investir dans leur propre infrastructure. La capitalisation boursière des grandes entreprises technologiques a augmenté depuis le début de la pandémie de COVID-19, influençant leur accumulation et leur concentration d'infrastructures informatiques, de données et de ressources humaines nécessaires à la formation de systèmes d'IA avancés (755). Cette accumulation de ressources est motivée par la réévaluation par les entreprises des rendements attendus des investissements dans l'IA.

Un système d'IA unique peut être adopté dans des secteurs critiques tels que la finance, la santé et la cybersécurité, ce qui accentue particulièrement les risques systémiques associés à la concentration du marché. Ces secteurs, qui sont interdépendants et essentiels à la sécurité nationale et à la stabilité économique, s'appuient de plus en plus sur l'IA pour la prise de décision, la détection des menaces, l'automatisation et l'optimisation des ressources. Les modèles d'IA à usage général dominants fournis par quelques grandes entreprises sont utilisés comme épine dorsale pour bon nombre de ces applications, ce qui crée un potentiel de vulnérabilités importantes (757). L'une des principales préoccupations est que les défauts, les vulnérabilités, les bugs ou les biais inhérents à ces systèmes d'IA largement adoptés pourraient entraîner des défaillances simultanées dans plusieurs secteurs (758).

Différents scénarios ont été proposés pour illustrer les perturbations potentielles. Par exemple, une faille de cybersécurité dans un modèle d'IA dominant pourrait exposer plusieurs institutions financières, agences gouvernementales et autres systèmes critiques à des cyberattaques coordonnées ou à des pannes de système (759, 760).

L'intensification du développement de normes techniques pour identifier et atténuer les points de défaillance uniques dans l'IA pourrait réduire les risques. Une façon d'atténuer les risques liés aux points de défaillance uniques consiste à réduire la probabilité que des modèles individuels échouent ou soient dangereux d'une manière ou d'une autre. Parmi les exemples de mesures d'atténuation potentielles explorées par les chercheurs figurent l'élaboration de normes techniques (761), ainsi que les exigences en matière d'audit et de reporting (762). Cependant, ces mesures d'atténuation impliquent des coûts et une complexité importants (763). Pour une discussion plus détaillée de diverses techniques de ce type, voir [3. Approches techniques de la gestion des risques.](#)

L'absence de méthodes établies pour modéliser les impacts à partir de points de défaillance uniques dans l'IA constitue un manque majeur de données probantes sur les risques de concentration du marché. Il est donc difficile de développer des méthodes d'atténuation fiables. Il est difficile de prédire comment les défaillances se propagent dans des systèmes sociétaux complexes. Il est donc difficile de prévoir de manière fiable les perturbations potentielles ou d'en comprendre toute l'ampleur. Cette incertitude entrave les efforts visant à concevoir des mesures de protection ciblées, car les données complètes destinées aux décideurs politiques et aux développeurs sur l'emplacement des vulnérabilités et la manière dont elles se manifestent sont encore en cours d'élaboration (764). Par conséquent, il existe un risque de stratégies d'atténuation incomplètes ou inefficaces, ce qui expose les secteurs critiques à un risque continu de défaillances en cas

Défauts du système. Bien que les chercheurs aient commencé à développer des méthodes pour mesurer la fiabilité des systèmes d'IA (765), elles sont peu nombreuses et leur adoption est limitée.

L'un des principaux défis auxquels sont confrontés les décideurs politiques qui cherchent à réduire les risques liés à la concentration du marché de l'IA à usage général est que le développement de cette technologie nécessite beaucoup de capitaux, ce qui favorise la domination de quelques très grands acteurs. La dynamique courante qui suit les tentatives de réduction de la concentration du marché illustre ce phénomène, les petites entreprises devenant rapidement des cibles d'acquisition pour des concurrents plus importants. Par exemple, le financement et les ressources peuvent aider les petites entreprises à se développer, mais cela tend à les rendre en retour des opportunités d'acquisition intéressantes pour les entreprises technologiques dominantes qui cherchent à éliminer la concurrence ou à étendre leurs capacités en matière d'IA (766, 767).

2.3.4. Risques pour l'environnement

INFORMATIONS CLÉS

- L'IA à usage général contribue modérément, mais rapidement, aux impacts environnementaux mondiaux par le biais de la consommation d'énergie et des émissions de gaz à effet de serre (GES). Les estimations actuelles indiquent que les centres de données et la transmission de données représentent environ 1 % des émissions mondiales de GES liées à l'énergie, l'IA consommant 10 à 28 % de la capacité énergétique des centres de données.
La demande énergétique de l'IA devrait augmenter considérablement d'ici 2026, certaines estimations prévoyant un doublement ou plus, principalement en raison des systèmes d'IA à usage général tels que les modèles linguistiques.
- Les avancées récentes en matière de capacités d'IA à usage général sont en grande partie dues à une augmentation marquée de la quantité de calcul nécessaire au développement et à l'utilisation de modèles d'IA, qui consomment davantage d'énergie. Alors que les entreprises d'IA alimentent de plus en plus leurs centres de données avec des énergies renouvelables, une part importante de la formation à l'IA dans le monde repose encore sur des sources d'énergie à forte teneur en carbone telles que le charbon ou le gaz naturel, ce qui entraîne les émissions susmentionnées et contribue au changement climatique.
- Le développement et le déploiement de l'IA ont également des impacts environnementaux importants, en raison de la consommation d'eau et de ressources, et des applications d'IA qui peuvent soit nuire, soit favoriser les efforts de développement durable. L'IA consomme de grandes quantités d'eau pour la production d'énergie, la fabrication de matériel et le refroidissement des centres de données. Toutes ces exigences augmentent proportionnellement au développement, à l'utilisation et aux capacités de l'IA. L'IA peut également être utilisée pour faciliter des activités préjudiciables à l'environnement, telles que l'exploration pétrolière, ainsi que dans des applications respectueuses de l'environnement susceptibles d'atténuer le changement climatique ou d'aider la société à s'y adapter, comme l'optimisation des systèmes de production et de transmission d'énergie.
- Les mesures d'atténuation actuelles comprennent l'amélioration du matériel, des logiciels et de l'énergie algorithmique. L'efficacité énergétique et le passage à des sources d'énergie sans carbone ont été des stratégies insuffisantes pour réduire les émissions de GES. L'amélioration de l'efficacité technologique et l'adoption des énergies renouvelables n'ont pas suivi le rythme de la demande énergétique : les émissions de GES des entreprises technologiques augmentent souvent malgré les efforts considérables déployés pour atteindre les objectifs de zéro émission nette de carbone. Des avancées technologiques importantes dans le matériel ou les algorithmes d'IA à usage général, ou des changements substantiels dans la production, le stockage et la transmission d'électricité, seront nécessaires pour répondre à la demande future sans que les impacts environnementaux n'augmentent en même temps.
- Depuis la publication du rapport intermédiaire (mai 2024), des preuves supplémentaires montrent que la demande d'énergie pour alimenter les charges de travail de l'IA augmente considérablement. Les développeurs d'IA à usage général ont signalé de nouveaux défis pour respecter leurs engagements en matière de carbone zéro net en raison de la consommation d'énergie accrue résultant du développement et de la fourniture de modèles d'IA à usage général, certains signalant une augmentation des émissions de GES en 2023 par rapport à 2022. En réponse, certaines entreprises se tournent vers l'énergie nucléaire pratiquement sans carbone pour alimenter les données de l'IA centres.

- Les principales lacunes en matière de données probantes concernant la consommation d'énergie et les émissions de GES de l'IA à usage général sont les suivantes : l'absence d'estimations précises de la consommation totale d'énergie ou des émissions dues à l'IA à usage général, et la difficulté d'anticiper les tendances futures correspondantes. Il n'existe pas suffisamment d'informations publiques sur les tendances actuelles de la consommation d'énergie de l'IA, comme la quantité de capacité des centres de données pouvant être attribuée à l'IA à usage général par rapport à d'autres charges de travail, et la quantité d'énergie ou d'autres impacts environnementaux pouvant être attribués à différents cas d'utilisation ou capacités de l'IA. Les chiffres actuels reposent en grande partie sur des estimations, qui deviennent encore plus variables et peu fiables lorsqu'elles sont extrapolées dans le futur en raison du rythme rapide du développement dans le domaine.

Définitions clés

- Émissions de GES (gaz à effet de serre) : rejet de gaz tels que le dioxyde de carbone (CO₂), le méthane, l'oxyde nitreux et les hydrofluorocarbures qui créent une barrière emprisonnant la chaleur dans l'atmosphère. Un indicateur clé du changement climatique.
- Intensité carbone : quantité d'émissions de GES produites par unité d'énergie. Utilisée pour quantifier les émissions relatives des différentes sources d'énergie.
- Calcul : abréviation de « ressources informatiques », qui fait référence au matériel (par exemple les GPU), aux logiciels (par exemple les logiciels de gestion des données) et à l'infrastructure (par exemple les centres de données) nécessaires pour former et exécuter les systèmes d'IA.
- Centre de données : un grand ensemble de serveurs informatiques en réseau à haute puissance utilisés pour le calcul à distance. Les centres de données hyperscale contiennent généralement plus de 5 000 serveurs.
- Effet rebond : en économie, la réduction des améliorations attendues en raison d'une augmentation de l'efficacité, résultant de changements corrélés dans le comportement, les modes d'utilisation ou d'autres changements systémiques. Par exemple, une amélioration de l'efficacité du moteur à combustion automobile (km/litre) de 25 % entraînera une réduction de moins de 25 % des émissions, car la réduction correspondante du coût de l'essence par kilomètre parcouru rendra la conduite plus économique, ce qui limitera les améliorations.
- Compensation carbone : compenser les émissions de GES d'une source en investissant dans d'autres activités qui empêchent des quantités comparables d'émissions ou éliminent le carbone de l'atmosphère, comme l'expansion des forêts.
- Transparence institutionnelle : la mesure dans laquelle les entreprises d'IA divulguent des informations techniques ou des informations organisationnelles soumises à un contrôle public ou gouvernemental, y compris des données de formation, des architectures de modèles, des données sur les émissions, des mesures de sécurité et de sûreté ou des processus décisionnels.

Les progrès récents dans les capacités d'IA à usage général ont été en grande partie alimentés par une augmentation rapide de la quantité de calcul nécessaire au développement et à l'utilisation des modèles d'IA. La méthodologie la plus simple pour améliorer les performances de l'IA à usage général sur les tâches finales consiste à permettre au modèle d'apprendre à partir d'autant d'exemples de données que possible. Cela est réalisé en augmentant la taille du modèle, mesurée en nombre de paramètres, à peu près proportionnellement à la quantité de données disponibles (156*, 157*). Pour qu'un modèle plus grand puisse apprendre ses paramètres à partir des données (en formation et en développement) et utiliser ces paramètres pour produire des résultats sur de nouvelles données

(en cours de déploiement ou d'utilisation), il doit effectuer davantage de calculs, ce qui nécessite davantage de puissance de calcul (voir [1.3. Capacités dans les années à venir](#) pour une discussion plus approfondie).

L'IA à usage général nécessite une énergie considérable pour se développer et s'utiliser, avec les émissions de GES et les impacts correspondants sur le réseau électrique. Par exemple, Meta estime que l'énergie nécessaire à l'entraînement de sa récente famille de LLM Llama 3 (juillet 2024) a entraîné 11 380 tonnes d'équivalent CO₂ (tCO₂e) d'émissions sur les quatre modèles publiés (11*). Les émissions totales correspondent à l'énergie consommée par 1 484 foyers américains moyens pendant un an, ou 2 708 véhicules de tourisme à essence conduits pendant un an (768). Google rapporte que l'entraînement de sa famille de LLM open source Gemma 2 a émis 1 247,61 tCO₂e (769*), mais comme la plupart des développeurs d'IA à usage général, ils ne divulguent pas la quantité d'énergie ou d'émissions nécessaires pour alimenter les modèles de production. Une énergie supplémentaire est nécessaire pour alimenter les centres de données dans lesquels la plupart des calculs d'IA à usage général sont effectués, notamment pour le refroidissement. Français Cette surcharge énergétique supplémentaire est généralement quantifiée sous la forme d'efficacité énergétique (PUE), qui est un ratio entre la quantité d'énergie utilisée pour le calcul et pour d'autres utilisations au sein d'un centre de données ; le PUE théorique optimal, indiquant une surcharge énergétique nulle, est de 1,0. Les centres de données hyperscale les plus efficaces, y compris de nombreux centres de données alimentant l'IA à usage général, affichent actuellement un PUE d'environ 1,1, la moyenne du secteur oscillant autour de 1,6 (770). La consommation d'énergie résulte également de la transmission de données sur les réseaux informatiques, qui est nécessaire pour communiquer les entrées et les sorties des modèles d'IA entre les appareils des utilisateurs, tels que les ordinateurs portables et les téléphones mobiles, et les centres de données où les modèles d'IA sont exécutés. Environ 260 à 360 TWh d'énergie étaient nécessaires pour prendre en charge les réseaux mondiaux de transmission de données en 2022, une quantité similaire à celle utilisée pour alimenter les centres de données (240 à 340 TWh, hors extraction de cryptomonnaies, qui représentait 100 à 150 TWh supplémentaires) la même année (771). À eux seuls, Google, Meta, Amazon et Microsoft, leaders dans la fourniture d'IA à usage général et d'autres services de cloud computing, étaient collectivement responsables de 69 % du trafic mondial de transmission de données, ce qui représente un changement par rapport aux années précédentes, où la majorité des transmissions de données étaient attribuées aux fournisseurs de services Internet publics (772).

Français Bien que les rapports se concentrent souvent sur le coût énergétique de la formation des modèles, il existe des preuves solides qu'une demande énergétique plus élevée découle de leur utilisation quotidienne. La formation et le développement correspondent à un nombre inférieur d'activités à forte consommation d'énergie, tandis que le déploiement correspond à un nombre très élevé d'utilisations à faible consommation d'énergie (puisque chaque requête utilisateur représente un coût énergétique) (739, 773, 774). Alors que les estimations les plus fiables de la consommation d'énergie et des émissions de GES dues à l'IA à usage général mesurent généralement leurs coûts de formation, comme ceux cités ci-dessus, les rapports disponibles suggèrent une proportion globale plus élevée de dépenses énergétiques dues à l'utilisation. En 2022, Google et Meta ont signalé que l'utilisation de systèmes d'IA représentait 60 à 70 % de l'énergie associée à leurs charges de travail d'IA, contre 0 à 40 % pour la formation et 10 % pour le développement (c'est-à-dire la recherche et l'expérimentation) (199, 206).

Le prétraitement et la génération de données pour l'IA à usage général entraînent également des coûts énergétiques importants. Meta a en outre signalé que le traitement des données, c'est-à-dire le filtrage et la conversion des données aux formats appropriés pour la formation des modèles d'IA, représentait 30 % de l'empreinte énergétique d'un modèle de production développé

Français en 2021 pour la recommandation et le classement personnalisés, et le calcul global consacré au prétraitement des données a augmenté de 3,2 fois entre 2019 et 2021 (199). Les grands modèles d'IA à usage général donnent lieu à plus de calculs pour le traitement des données que les modèles d'IA étroits. Non seulement les modèles d'IA à usage général consomment beaucoup plus de données que les modèles étroits, mais les modèles eux-mêmes sont de plus en plus utilisés pour générer des données synthétiques supplémentaires pendant le processus de formation et pour choisir les meilleures données synthétiques sur lesquelles s'entraîner (37*, 775, 776*). Ils sont également utilisés pour générer des données pour la formation de modèles d'IA étroits (777). Cependant, des chiffres récents fournissant une attribution aussi détaillée de la consommation d'énergie de l'IA à usage général ne sont pas disponibles. La disponibilité limitée de données plus larges quantifiant la consommation d'énergie de l'IA a donné lieu à des mandats récents, comme dans la loi européenne sur l'IA, se concentrant sur la formation des modèles malgré la nécessité d'un reporting et d'une caractérisation accrue des exigences dues au traitement des données et à l'utilisation des modèles (778).

Actuellement, les émissions de GES de l'IA à usage général proviennent principalement de l'intensité carbone des sources d'énergie utilisées pour alimenter les centres de données et les réseaux de transmission de données qui soutiennent leur formation et leur utilisation. Par exemple, les sources renouvelables telles que l'énergie solaire émettent beaucoup moins de GES que les combustibles fossiles (779*). Alors que les entreprises d'IA alimentent de plus en plus leurs centres de données avec de l'énergie renouvelable (199, 206, 780*, 781), une part importante des calculs d'IA à l'échelle mondiale repose encore sur des sources à forte teneur en carbone telles que le charbon ou le gaz naturel (779*). Cela entraîne d'importantes émissions de GES.

Les estimations de la consommation totale d'énergie et des émissions de GES liées aux centres de données et à l'IA varient. Selon les estimations de l'Agence internationale de l'énergie (AIE), les centres de données et la transmission de données représentent 1 % des émissions mondiales de GES liées à la consommation d'énergie et 0,6 % de toutes les émissions de GES (qui incluent également d'autres sources de GES telles que l'agriculture et les processus industriels). (770, 771, 782). Selon des estimations récentes, entre 10 et 28 % de la consommation énergétique des centres de données est due à l'utilisation de l'IA, principalement en raison de l'IA générative (LLM et modèles de génération d'images), qui représente la majeure partie de la consommation énergétique due à l'IA à usage général (770, 771, 782). La combinaison de ces estimations suggère que l'utilisation de l'IA est responsable de 0,1 à 0,28 % des émissions mondiales de GES attribuées à la consommation d'énergie et de 0,06 à 0,17 % de toutes les émissions de GES, mais les pourcentages exacts dépendent de la quantité d'énergie utilisée provenant de sources d'énergie à forte intensité de carbone. L'intensité carbone moyenne de l'électricité alimentant les centres de données aux États-Unis est de 548 grammes de CO₂ par kWh, ce qui est presque 50 % supérieur à la moyenne nationale américaine (783). Les facteurs affectant les émissions de GES comprennent l'emplacement des centres de données et l'heure de la journée de consommation d'énergie, l'efficacité du centre de données et l'efficacité du matériel utilisé. Par conséquent, les émissions réelles de GES pour une quantité donnée d'énergie consommée par l'IA peuvent varier considérablement.

Depuis la publication du rapport intermédiaire, de nouvelles preuves de l'augmentation de la demande énergétique pour alimenter les centres de données exécutant des charges de travail d'IA ont été observées. En octobre 2024, l'AIE prévoit que les centres de données représenteront moins de 10 % de la croissance de la demande mondiale d'électricité entre 2023 et 2030 (784). La majeure partie de la croissance globale de la demande devrait provenir d'autres sources croissantes de demande d'électricité, telles que l'adoption des véhicules électriques et les besoins accrus en matière de refroidissement des bâtiments. Cependant, les impacts sur les centres de données sont très localisés par rapport à d'autres secteurs, ce qui entraîne une répartition inégale de la demande accrue et

Français impacts disproportionnés dans certains domaines (784). Par exemple, les centres de données ont consommé plus de 20 % de toute l'électricité en Irlande en 2023 (785), et la consommation d'électricité augmente aux États-Unis, où se trouvent plus de la moitié de la capacité mondiale des centres de données (786), pour la première fois depuis plus d'une décennie, en partie en raison du développement et de l'utilisation accrue de l'IA (787). Les entreprises technologiques se tournent vers l'énergie nucléaire (qui présente ses propres avantages et risques complexes) comme source d'énergie neutre en carbone pour alimenter les centres de données, de nombreuses grandes entreprises technologiques signant des accords avec des fournisseurs d'électricité pour sécuriser l'énergie nucléaire. En septembre 2024, Microsoft a signé un accord qui rouvrira la centrale nucléaire de Three Mile Island en Pennsylvanie, s'engageant à acheter toute la capacité de production de la centrale pour les 20 prochaines années, soit suffisamment pour alimenter environ 800 000 foyers (788*). Amazon a signé un accord similaire en mars pour acheter jusqu'à 960 MW/an d'énergie nucléaire pour alimenter un campus de centres de données pour sa plateforme cloud Amazon Web Services (AWS) (789), ce qui représente le premier cas de colocalisation d'un centre de données avec une centrale nucléaire. Cependant, en novembre, la Federal Energy Regulatory Commission américaine a rejeté la demande du fournisseur de transmission de modifier son accord de service d'interconnexion pour augmenter la transmission vers le centre de données (790), jetant des doutes sur la question de savoir si les régulateurs soutiendront une telle colocalisation à l'avenir. En octobre, Google a annoncé un accord d'achat d'énergie nucléaire à partir de petits réacteurs modulaires (SMR), le premier accord d'entreprise de ce type au monde, affirmant qu'ils avaient besoin de cette nouvelle source d'électricité pour « soutenir les technologies d'IA » (791*).

Il existe plusieurs mesures d'atténuation potentielles pour réduire la consommation d'énergie et les émissions de GES croissantes des systèmes d'IA à usage général, comme le passage à une énergie sans carbone, l'achat de compensations carbone et l'amélioration de l'efficacité des systèmes d'IA et des centres de données, mais il n'existe pas de solution miracle. Comme dans d'autres secteurs, continuer à transférer l'énergie des centres de données d'IA à usage général vers des sources d'énergie renouvelables telles que l'éolien, l'hydroélectrique et le solaire est une voie prometteuse, mais elle est actuellement limitée par la technologie de stockage et de transmission des batteries ; les sources renouvelables ne peuvent actuellement pas fournir d'énergie aux centres de données qui en ont besoin sans interruption, dans des régions géographiquement diverses. Comme mentionné dans le paragraphe précédent, les entreprises d'IA manifestent un intérêt accru pour les sources d'énergie nucléaire, en particulier les SMR moins chers et plus sûrs, pour combler le vide à court et moyen terme. Alors que les SMR fournissent une énergie sans carbone ininterrompue, un rapport récent souligne que les SMR (< 300 MW) produisent plus de déchets nucléaires par unité d'énergie produite que les réacteurs à grande échelle (> 1 000 MW) d'un facteur 2 à 30 (792). L'amélioration de l'efficacité énergétique des systèmes d'IA à usage général, mesurée en termes d'énergie utilisée pour atteindre un niveau de capacité donné, est un autre moyen de réduire la consommation d'énergie (206, 773). Une allocation et une planification plus intelligentes des ressources constituent également une voie prometteuse pour réduire les émissions de GES.

Les charges de travail d'IA à usage général peuvent être suspendues pendant les périodes de pointe de consommation d'énergie afin de réduire les émissions de GES de près de 30 % pour certaines régions et certains bouquets énergétiques (793), mais toutes les charges de travail d'IA à usage général ne sont pas compatibles avec cette approche, en particulier l'inférence de modèle, qui nécessite généralement que la charge de travail soit exécutée immédiatement afin de renvoyer une réponse à l'utilisateur immédiatement (par exemple lorsqu'un résumé d'IA à usage général est inclus dans le cadre d'une recherche sur le Web). D'autres stratégies d'atténuation comprennent l'élaboration d'évaluations d'impact sur la durabilité pour le développement et le déploiement de l'IA ; des restrictions ou des conditions de ressources pour la formation de l'IA ; et des budgets énergétiques négociables pour la formation et l'inférence de l'IA (794).

Les compensations carbone sont une méthode populaire utilisée par les développeurs d'IA à usage général pour atténuer les émissions de GES, mais elles n'entraînent pas toujours de réelles réductions d'émissions. Les consommateurs d'énergie tentent généralement d'atténuer leurs émissions de GES en concluant des accords d'achat d'énergie renouvelable (PPA), des crédits d'énergie renouvelable (REC), des mécanismes de transition vers le charbon (CTM) ou des certificats de compensation carbone afin de compenser leurs émissions en achetant une énergie renouvelable équivalente, ou en investissant dans d'autres projets de réduction des émissions de carbone ou de transition vers l'énergie verte. Les compensations carbone sont le principal mécanisme actuellement utilisé par les entreprises technologiques pour atteindre leurs engagements de zéro émission nette de carbone, parallèlement à l'augmentation de l'approvisionnement en sources d'énergie renouvelables pour alimenter directement la consommation énergétique des centres de données (780*, 795*, 796*). Par exemple, Meta rapporte qu'ils ont atténué les émissions grâce à la formation LLM citée ci-dessus en achetant une quantité équivalente d'énergie renouvelable (11*). Cette stratégie présente également des limites, en raison de la difficulté de vérifier l'additionnalité des projets de compensation, c'est-à-dire de garantir que les réductions d'émissions n'auraient pas eu lieu quel que soit le programme de compensation (797).

Français L'efficacité énergétique de l'IA à usage général s'améliore rapidement, mais pas suffisamment pour arrêter la croissance continue des émissions. Le matériel d'IA spécialisé et d'autres améliorations de l'efficacité matérielle améliorent les performances par watt des charges de travail d'apprentissage automatique au fil du temps (206). De plus, les nouvelles techniques et architectures d'apprentissage automatique peuvent également contribuer à réduire la consommation d'énergie (206), tout comme les améliorations des cadres logiciels et des algorithmes de support (798, 799). L'énergie utilisée par unité de calcul a été réduite d'environ 26 % par an (144). Cependant, les taux actuels d'amélioration de l'efficacité sont insuffisants pour répondre à la demande croissante. La demande de puissance de calcul utilisée pour la formation de l'IA, qui a augmenté d'un facteur d'environ 4 fois chaque année, dépasse jusqu'à présent largement les améliorations de l'efficacité énergétique (26). Cette inadéquation se reflète dans le fait que les entreprises technologiques impliquées dans le développement et le déploiement de l'IA à usage général signalent des difficultés à atteindre les objectifs de durabilité environnementale. Baidu rapporte que l'augmentation des besoins énergétiques due au « développement rapide des LLM » pose de « sérieux défis » au développement de centres de données écologiques (781), et Google signale également une augmentation de 17 % de la consommation énergétique des centres de données en 2023 par rapport à 2022 et une augmentation de 37 % des émissions de GES dues à la consommation d'énergie « malgré des efforts et des progrès considérables en matière d'énergie sans carbone ». Ils attribuent ces augmentations à l'augmentation des investissements dans l'IA (780*).

Les améliorations de l'efficacité énergétique ne suffisent pas à elles seules à annuler la croissance globale de la consommation énergétique de l'IA et pourraient même l'accélérer davantage en raison des « effets de rebond ». Les économistes ont constaté que pour les technologies précédentes, les améliorations de l'efficacité énergétique tendent à augmenter, plutôt qu'à diminuer, la consommation globale d'énergie en diminuant le coût par unité de travail (800). Les améliorations de l'efficacité peuvent conduire à une plus grande consommation d'énergie en rendant des technologies telles que l'IA à usage général moins chères et plus facilement disponibles, et en augmentant la croissance du secteur.

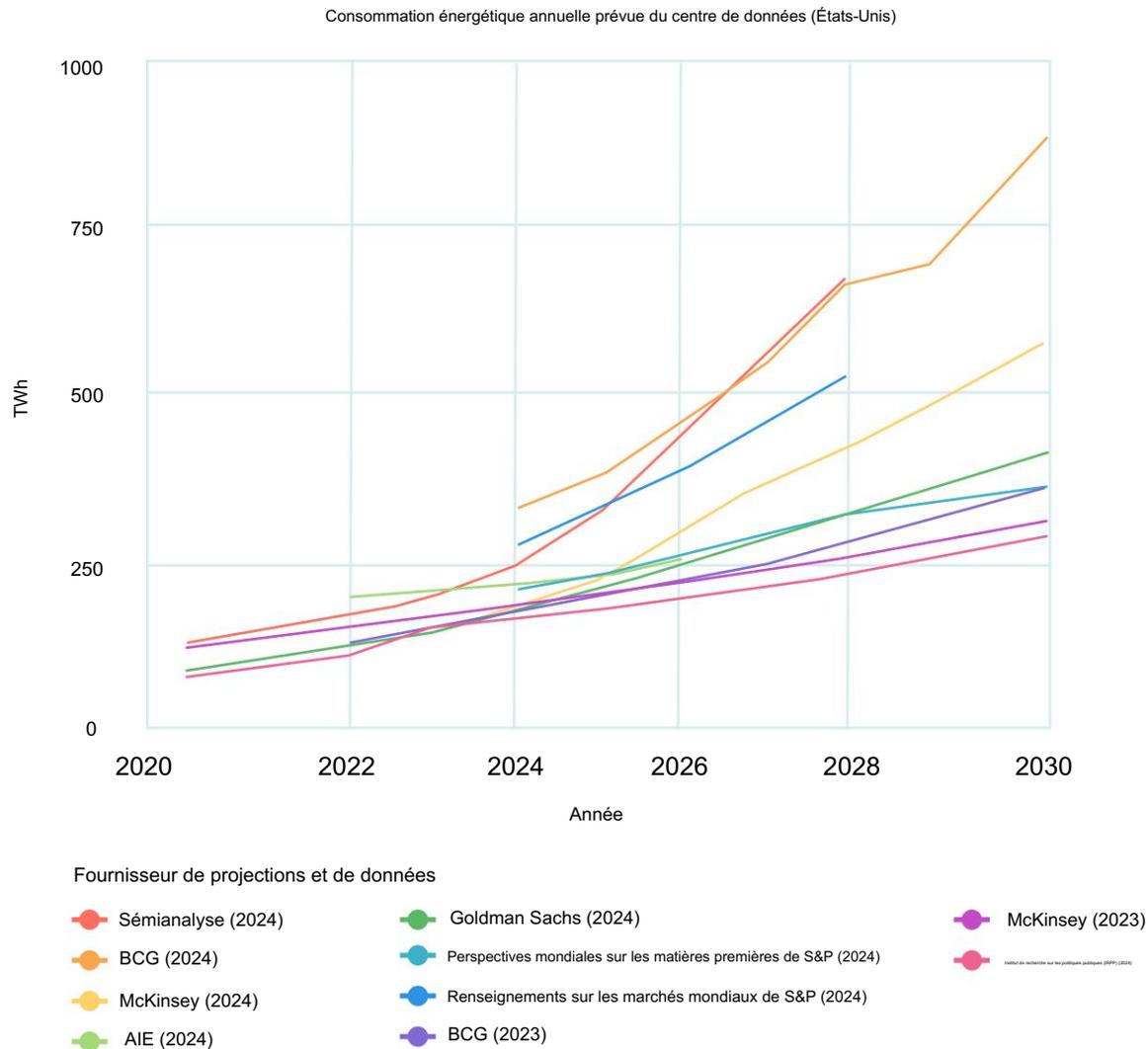


Figure 2.10 : La consommation énergétique des centres de données aux États-Unis devrait croître rapidement, atteignant entre 270 et 930 TWh par an d'ici 2030. Cette large gamme de projections (variant de plus de 600 TWh, soit plus de 10 % de la consommation totale d'énergie des États-Unis en 2022) résulte de l'évolution rapide de la technologie et de données historiques limitées, en particulier pour l'utilisation spécifique de l'IA. Source : Kamiya, G. & Coroamã, VC, 2024 (801).

Les principales lacunes en matière de données sur la consommation d'énergie et les émissions de l'IA à usage général sont le manque d'estimations précises de la consommation totale d'énergie, des émissions ou de la consommation de ressources dues à l'IA à usage général, et la difficulté d'anticiper les tendances futures correspondantes. Les estimations ascendantes de la consommation d'énergie et des émissions, telles que celles décrites ci-dessus, sont beaucoup plus faciles à calculer que les estimations descendantes pour l'ensemble du secteur. On pense généralement que le développement et l'utilisation accrues de l'IA à usage général entraînent une augmentation de la demande de capacité de calcul des centres de données et de la consommation d'énergie correspondante, et on suppose donc que les tendances globales de la consommation d'énergie des centres de données reflètent la croissance du développement et de l'utilisation de l'IA à usage général. En 2023, les centres de données (à l'exclusion de l'extraction de cryptomonnaies) représentaient entre 1 % et 1,5 % de la demande mondiale d'électricité (802) ; environ 2 % dans l'UE, 4 % aux États-Unis et près de 3 % en Chine (213, 803, 804). En 2020, les centres de données et les réseaux de transmission de données ont émis 330 millions

tCO₂e au total, soit un peu moins de 1 % de toutes les émissions de GES liées à l'énergie et 0,6 % des émissions mondiales de GES (771). Alors que l'on estime actuellement que l'IA représente 10 à 30 % des charges de travail des centres de données (770, 782), la demande de développement et d'utilisation de l'IA (à usage général et autre) devrait continuer de croître dans les années à venir. Certaines sources estiment que cette croissance doublera la demande d'électricité des centres de données, passant de 460 TWh en 2022 à plus de 1 000 TWh en 2026 (208).

Français Google rapporte avoir émis 14,3 millions de tCO₂e en 2023, soit une augmentation de 13 % par rapport à 2022 et de 48 % depuis 2019, qu'elle attribue à l'augmentation de la consommation énergétique des centres de données liée à l'intégration accrue de l'IA dans les produits (780*). Cependant, les projections varient considérablement et il est fondamentalement difficile de projeter l'utilisation et la croissance futures de l'IA en raison du rythme rapide et imprévisible du développement de la technologie (805). La figure 2.10 illustre la large gamme d'estimations disponibles pour la future consommation énergétique des centres de données aux États-Unis, qui varient considérablement, jusqu'à 10 % de la consommation totale d'électricité aux États-Unis en 2022. Il est typique d'estimer les tendances futures en extrapolant simplement à partir de la demande actuelle et du taux de croissance d'un indicateur. Cependant, cette méthodologie ignore les variables critiques qui dictent la croissance réelle et s'est avérée insuffisante pour estimer avec précision la demande en raison des développements technologiques. Par exemple, alors que le trafic Internet mondial (un indicateur de la consommation électrique des centres de données) a été multiplié par plus de dix entre 2010 et 2020, la consommation électrique des centres de données n'a augmenté que de 6 % sur la même période en raison des améliorations apportées au matériel et à l'efficacité des centres de données.

Français Une percée technologique dans les algorithmes d'IA à usage général d'aujourd'hui pourrait également réduire les besoins énergétiques, et les estimations actuelles de la croissance de la consommation d'énergie doivent prendre en compte des facteurs supplémentaires qui freinent la croissance, tels que les limitations de la chaîne d'approvisionnement en matériel d'IA et la capacité de production d'électricité (784). En outre, l'IA (à usage général ou non) peut également avoir des impacts environnementaux indirects (positifs ou négatifs) découlant d'applications spécifiques (774). Par exemple, l'IA pourrait être appliquée pour accélérer la découverte en science des matériaux d'une nouvelle chimie de batterie permettant une adoption plus large des énergies renouvelables, ou pour identifier des catalyseurs permettant une capture plus efficace du carbone ou une production d'hydrogène combustible (807). L'IA peut également être appliquée à des objectifs préjudiciables à l'environnement tels que l'exploration et l'extraction de pétrole et de gaz, entraînant une augmentation des émissions de GES (807). La quantification des impacts indirects est encore plus difficile que la caractérisation de ses impacts directs en raison, par exemple, de la consommation d'énergie, et des travaux supplémentaires sont nécessaires pour développer des cadres solides pour l'évaluation du cycle de vie des modèles d'IA à usage général (774). Il est nécessaire de mieux rendre compte et caractériser les demandes énergétiques passées et actuelles, ainsi que les cas d'utilisation dominants de l'IA qui les alimentent, afin d'évaluer les risques et d'élaborer des stratégies d'atténuation pour l'augmentation de la demande énergétique et des émissions dues à l'IA à usage général (778). Afin de rester sur la bonne voie avec le scénario de zéro émission nette d'ici 2050 de l'AIE, par exemple, les émissions dues aux centres de données et à la transmission de données devraient être réduites de moitié d'ici 2030 (771), mais on ne sait pas bien quelle proportion de ces émissions peut être attribuée à l'IA à usage général, quels développements et cas d'utilisation de l'IA à usage général contribuent le plus à ces émissions et lesquels atténuent ou réduisent les émissions ailleurs, et comment ces tendances évoluent au fil du temps.

Outre les émissions de GES dues à la consommation d'énergie, l'IA à usage général a d'autres impacts environnementaux dus aux systèmes et structures physiques nécessaires à son développement et à son utilisation, qui sont encore moins bien compris. Les émissions de GES dues à la consommation d'énergie évoquées jusqu'à présent sont généralement appelées émissions opérationnelles, et elles représentent actuellement la plus grande proportion d'émissions. L'empreinte carbone incorporée du matériel d'IA, qui comprend les émissions provenant de la fabrication, du transport, de l'infrastructure physique du bâtiment et de l'élimination, contribue également à des émissions de GES importantes. Selon l'emplacement et le scénario, cela peut représenter jusqu'à 50 % des émissions totales d'un modèle (199). À mesure que l'efficacité énergétique opérationnelle s'améliore, l'empreinte carbone incorporée deviendra une proportion plus importante de l'empreinte carbone totale (808). Intel indique que son campus d'Ocotillo a généré plus de 200 000 tCO₂e en 2023 à partir des seules émissions directes (hors électricité) (809*), et est en passe de générer plus de 300 000 tCO₂e d'ici la fin de 2024, après avoir consommé plus d'un milliard de kWh d'énergie au premier trimestre 2024 (809*). L'estimation de l'empreinte carbone incorporée actuelle de l'IA à usage général pose un grand défi en raison du manque de données des fabricants de matériel. Cela résulte d'une combinaison d'incitations, notamment la volonté des fabricants de protéger la propriété intellectuelle autour des processus de fabrication propriétaires et la consolidation de l'expertise dans la fabrication de matériel d'IA spécialisé dans un nombre très limité d'entreprises, ce qui limite l'accès et le transfert des connaissances.

La consommation d'eau est un autre domaine émergent de risque environnemental lié à l'IA à usage général.

Le développement et l'utilisation de l'IA à usage général prélèvent de l'eau douce dans les réseaux d'eau locaux, dont une partie est ensuite consommée, principalement par évaporation. Comme pour la consommation d'énergie, la consommation d'eau de l'IA à usage général augmente également à mesure que les modèles deviennent plus grands. L'IA à usage général a des besoins en eau à la fois incorporés et opérationnels. La consommation d'eau incorporée provient de l'utilisation de l'eau dans le processus de fabrication du matériel, et la consommation d'eau opérationnelle provient principalement de la production d'énergie et des systèmes de refroidissement par évaporation dans les centres de données. Dans la production d'énergie, l'eau s'évapore lorsqu'elle est utilisée pour le refroidissement des centrales nucléaires et à combustion de combustibles fossiles et des barrages hydroélectriques. Dans les centres de données, le matériel informatique produit également une quantité importante d'énergie résiduelle sous forme de chaleur et doit être refroidi afin d'optimiser l'efficacité et la longévité des calculs.

Les méthodes les plus efficaces et les plus répandues pour refroidir le matériel dans les centres de données utilisent l'évaporation de l'eau.

À mesure que les calculs utilisés pour la formation et le déploiement de modèles d'IA à usage général augmentent, les besoins en refroidissement augmentent également, ce qui entraîne une augmentation de la consommation d'eau. L'eau est également consommée lors des processus de fabrication du matériel. En 2023, l'usine de fabrication de puces Intel d'Ocotillo en Arizona, qui a obtenu la plus haute certification en matière de conservation de l'eau de l'Alliance for Water Stewardship, a prélevé 10 561 millions de litres d'eau (90 % d'eau douce), dont 1 896 millions de litres ont été consommés (809*). En supposant une consommation moyenne d'eau par ménage de 144 litres par jour (810), cela équivaut au prélèvement d'eau annuel de plus de 200 000 ménages. Taiwan Semiconductor Manufacturing Company (TSMC), le plus grand fabricant mondial de semi-conducteurs et le principal fournisseur de puces pour les entreprises de matériel d'IA telles que Nvidia, rapporte qu'en 2023, sa consommation d'eau par unité avait augmenté de 25,2 % depuis 2010, malgré son objectif de réduire la consommation de 2,7 % sur cette période et de 30 % d'ici 2030 ; et ce malgré des mesures d'économie d'eau accrues qui ont permis à TSMC d'économiser 33 % d'eau de plus d'une année sur l'autre en 2023 (811). La consommation d'eau selon les modèles actuels et la méthodologie pour l'évaluer font encore l'objet d'un débat scientifique, mais certains

Les chercheurs prévoient que la consommation d'eau par l'IA pourrait atteindre des milliers de milliards de litres d'ici 2027 (199, 812). Dans le contexte des préoccupations concernant la pénurie mondiale d'eau douce, et sans avancées technologiques permettant des alternatives efficaces en termes d'émissions, l'empreinte hydrique de l'IA pourrait constituer une menace substantielle pour l'environnement et le droit humain à l'eau (813). En réponse aux mandats du Congrès, le ministère américain de l'Énergie travaille actuellement à l'évaluation des besoins actuels et futurs en consommation d'énergie et d'eau des centres de données, un rapport devant être publié d'ici la fin de 2024 (787). Les opérateurs de centres de données européens doivent déclarer leur consommation d'eau à partir de 2025 (814).

Les mesures d'atténuation potentielles de la consommation d'eau liée à l'IA comprennent la réduction de la consommation d'énergie et le développement et le déploiement de processus à faible consommation d'eau pour le refroidissement et la fabrication. Les mêmes améliorations algorithmiques et logicielles déployées pour atténuer la consommation d'énergie entraîneront également une certaine réduction de la consommation d'eau, car une partie de la consommation d'eau est due à la consommation d'énergie. D'autres mesures d'atténuation de la consommation d'énergie, telles que l'amélioration de l'efficacité du matériel ou le passage à des sources d'énergie sans carbone, n'entraîneront pas nécessairement une réduction de la consommation d'eau, et pourraient l'augmenter ; l'amélioration de l'efficacité du matériel implique la fabrication de nouveau matériel pour remplacer l'ancien, et la production d'énergie nucléaire nécessite plus d'eau pour le refroidissement que la production de gaz naturel (815). Des technologies plus récentes, telles que le refroidissement à sec, peuvent être utilisées pour réduire les prélèvements d'eau nécessaires au refroidissement des centrales électriques, mais le refroidissement à sec diminue l'efficacité de la production d'énergie (816). Dans les centres de données et la fabrication de matériel, l'eau peut être récupérée et recyclée, mais cela nécessite également un apport énergétique accru afin de filtrer l'eau jusqu'à une grande pureté, par exemple par osmose inverse (809*, 817). Ces exemples mettent en évidence un compromis courant entre la consommation d'énergie et les préoccupations relatives à la consommation d'eau qui doivent être prises en compte lors de l'élaboration de politiques concernant les impacts environnementaux de l'IA. Français Les centres de données peuvent être construits dans des régions géographiques au climat froid propices au refroidissement naturel par air, mais les défis logistiques en termes d'énergie et de transmission de données, de construction et de conditions météorologiques extrêmes limitent la rentabilité de cette approche à grande échelle. La trigénération, dans laquelle la chaleur résiduelle de la production d'énergie est utilisée pour assurer le refroidissement, peut minimiser la consommation d'eau et d'énergie dans les centres de données (818). Cependant, les systèmes de trigénération actuels sont généralement alimentés par la combustion de combustibles fossiles et des recherches supplémentaires sont nécessaires pour développer des systèmes de trigénération alimentés par des sources d'énergie sans carbone et à faible consommation d'eau. Le refroidissement par plasma d'hydrogène pourrait également améliorer l'efficacité du refroidissement des centres de données, mais des efforts importants sont encore nécessaires pour développer une infrastructure robuste pour la production d'hydrogène qui ne dépend pas des combustibles fossiles (819). En conjonction avec les efforts visant à optimiser les processus de fabrication, les fabricants de matériel informatique ont commencé à déclarer ou à s'engager à une utilisation d'eau « nette positive », en combinant la réduction de leur consommation d'eau et le financement de projets externes de restauration de l'eau équivalents en gallons à leur consommation, dans la même veine que les engagements de zéro émission nette de carbone qui tirent parti des REC ou des compensations carbone (809*, 811), avec des défis similaires.

Les décideurs politiques sont confrontés à trois défis majeurs pour faire face à l'impact de l'IA sur l'environnement : une transparence institutionnelle limitée sur les données relatives à la consommation d'énergie et aux émissions, des relations floues entre les coûts de calcul et la question de savoir si les capacités qui en résultent sont appliquées pour le bien ou le mal de l'environnement, et une grande incertitude due au développement rapide. Les données disponibles pour quantifier l'énergie et les émissions associées à l'IA à usage général sont limitées, ce qui limite la capacité des chercheurs à analyser et à prévoir les modèles d'utilisation et les politiques correspondantes.

Français développement. Cependant, les exigences de reporting existantes ne fournissent toujours pas suffisamment d'informations sur l'IA en particulier ; elles n'exigent pas des développeurs qu'ils décomposent les impacts par phases d'utilisation du modèle (formation versus utilisation) ou par cas d'utilisation (usage général versus spécifique à une tâche, ou si l'IA est appliquée pour atténuer ou accélérer les impacts environnementaux négatifs, comme pour aider à l'extraction de pétrole et de gaz) (778). En outre, la communauté scientifique ne comprend pas suffisamment la quantité de calcul nécessaire pour atteindre le niveau de capacité souhaité à partir d'un modèle d'IA à usage général. Cela limite la mesure dans laquelle des objectifs de consommation d'énergie peuvent être fixés pour des modèles ou des cas d'utilisation spécifiques, comme la quantité d'énergie ou d'émissions allouée pour générer une image, car les limites supérieure et inférieure de l'énergie requise sont soit très larges, soit très spécifiques à chaque cas. Une collaboration étroite et une communication efficace sont nécessaires entre les experts du domaine et les décideurs politiques pour garantir que les décisions politiques sont fondées sur des données précises et que des mécanismes sont mis en place pour garantir que de meilleures données soient disponibles à l'avenir pour soutenir l'élaboration et la mise en œuvre des pol

2.3.5. Risques pour la vie privée

INFORMATIONS CLÉS

- Les systèmes d'IA à usage général peuvent provoquer ou contribuer à des violations de la vie privée des utilisateurs.
Des violations peuvent survenir par inadvertance lors de la formation ou de l'utilisation de systèmes d'IA, par exemple en cas de traitement non autorisé de données personnelles ou de fuite de dossiers médicaux utilisés lors de la formation. Mais des violations peuvent également survenir délibérément par l'utilisation d'IA à usage général par des acteurs malveillants, par exemple s'ils utilisent l'IA pour déduire des faits privés ou violer la sécurité.
- L'IA à usage général divulgue parfois des informations sensibles acquises lors de l'entraînement ou de l'interaction avec les utilisateurs.
Les informations sensibles contenues dans les données d'entraînement peuvent être divulguées de manière involontaire lorsqu'un utilisateur interagit avec le modèle. De plus, lorsque les utilisateurs partagent des informations sensibles avec le modèle pour obtenir des réponses plus personnalisées, ces informations peuvent également être divulguées ou exposées à des tiers non autorisés.
- Les acteurs malveillants peuvent utiliser l'IA à usage général pour contribuer à la violation de la vie privée. Systèmes d'IA peut faciliter des recherches plus efficaces et efficaces de données sensibles et peut déduire et extraire des informations sur des individus spécifiques à partir de grandes quantités de données. Cette situation est encore exacerbée par les risques de cybersécurité créés par les systèmes d'IA à usage général (voir [2.1.3](#)).
[\(Cyber-infraction\)](#)
- Depuis la publication du rapport intermédiaire (mai 2024), les gens utilisent de plus en plus l'IA à usage général dans des contextes sensibles tels que les soins de santé ou la surveillance du lieu de travail. Cela crée de nouveaux risques pour la vie privée qui, jusqu'à présent, ne se sont toutefois pas matérialisés à grande échelle. En outre, les chercheurs tentent de supprimer les informations sensibles des données de formation et de créer des outils de déploiement sécurisés.
- Pour les décideurs politiques, il reste difficile de connaître l'ampleur ou la portée des violations de la vie privée.
Il est extrêmement difficile d'évaluer l'ampleur des atteintes à la vie privée causées par l'IA à usage général, car de nombreux préjudices surviennent de manière involontaire ou à l'insu des personnes concernées. Même dans le cas de fuites documentées, il peut être difficile d'identifier leur source, car les données sont souvent traitées sur plusieurs appareils ou dans différentes parties de la chaîne d'approvisionnement.

Définitions clés

- Confidentialité : droit d'une personne ou d'un groupe à contrôler la manière dont les autres accèdent à ses données sensibles ou les traitent. informations et activités.
- Informations personnelles identifiables (IPI) : toute donnée permettant d'identifier directement ou indirectement une personne (par exemple, un nom ou un numéro d'identification). Il s'agit d'informations pouvant être utilisées seules ou combinées à d'autres données pour identifier une personne de manière unique.
- Données sensibles : informations qui, si elles étaient divulguées ou mal traitées, pourraient entraîner des dommages, embarras, inconfort ou injustice envers un individu ou une organisation. • Minimisation des données : pratique consistant à collecter et à conserver uniquement les données directement nécessaires à un objectif spécifique, et à les supprimer une fois cet objectif atteint.

- Génération augmentée par récupération (RAG) : une technique qui permet aux LLM d'extraire des informations d'autres sources lors de l'inférence, telles que les résultats de recherche sur le Web ou une base de données interne de l'entreprise, permettant ainsi des réponses plus précises ou personnalisées.
- Deepfake : un type de faux contenu généré par l'IA, composé de contenu audio ou visuel, qui déforme les faits et présente des personnes réelles comme faisant ou disant quelque chose qu'elles n'ont pas réellement fait ou dit.

Les systèmes d'IA à usage général s'appuient sur de vastes quantités de données personnelles et peuvent les traiter, ce qui présente des risques importants pour la confidentialité. Dans le contexte de l'IA, la confidentialité est un concept complexe et multiforme qui englobe :

- Confidentialité des données et protection des données personnelles collectées ou utilisées à des fins de formation, réglage fin, extraction d'informations ou pendant l'inférence.
- Transparence institutionnelle et contrôles sur la manière dont les informations personnelles sont utilisées dans les systèmes d'IA (820) ; par exemple, la possibilité pour les individus de refuser que leurs données personnelles soient collectées à des fins de formation, ou la capacité a posteriori de faire en sorte qu'un système d'IA à usage général « désapprenne » des informations spécifiques sur un individu (821) ; et les défis connexes tels que la conciliation de la minimisation des données et de la transparence (822), le contrôle de la manière dont les décisions fondées sur les données sont prises, et l'utilisation ou le traitement non autorisé des données personnelles (823).
- Protection contre les préjudices individuels et collectifs pouvant résulter de l'utilisation ou de l'utilisation malveillante des données. Par exemple, la création de deepfakes (824), les atteintes au droit à l'oubli (548) ou au droit de rectification (825), et d'autres risques liés au scraping à grande échelle de données personnelles (826).

L'IA à usage général présente divers risques pour la vie privée. Ceux-ci sont classés de manière très générale en :

- Risques liés à la formation : risques liés à la formation et à la collecte de données (notamment de données sensibles).
- Risques d'utilisation : risques liés au traitement des informations sensibles par les systèmes d'IA lors de leur utilisation.
- Risques de préjudice intentionnel : risques que des acteurs malveillants utilisent l'IA à usage général pour nuire vie privée individuelle (voir figure 2.11).

Ces risques sont déjà présents avec les outils d'IA actuellement disponibles, mais sont exacerbés par l'ampleur accrue de la formation, la capacité de traitement de l'information et la facilité d'utilisation présentées par l'IA à usage général.

Les systèmes d'IA à usage général peuvent exposer leurs données de formation (« risques de formation »). La formation des modèles d'IA à usage général nécessite généralement de grandes quantités de données. Des études universitaires ont montré que certaines de ces données de formation peuvent être mémorisées par des modèles d'IA à usage général (827, 828), permettant aux utilisateurs de déduire des informations sur les personnes dont les données ont été collectées (829, 830, 831) ou même de reconstruire des exemples de formation entiers (832, 833, 834, 835). Cependant, les définitions de la mémorisation varient, il est donc difficile de faire des déclarations concrètes sur les préjudices qui pourraient découler de la mémorisation (827). De nombreux systèmes sont formés sur des données accessibles au public contenant des informations personnelles à l'insu ou sans le consentement des personnes concernées, en plus de la formation

sur le contenu Web propriétaire appartenant à des distributeurs de médias (826, 836). Cela s'étend aux cas où une personne publie des informations personnelles sur une autre personne en ligne – par exemple, des publications Facebook comprenant des photos et des informations sur les pairs ou les amis d'une personne sans le consentement explicite de ces pairs. Dans des domaines spécifiques, la formation sur des données sensibles (telles que des données médicales ou financières) est souvent nécessaire pour améliorer les performances dans ce domaine, mais pourrait entraîner de graves fuites de confidentialité.

Ces risques peuvent être réduits – par exemple, les systèmes d'IA à usage général existants, comme Gemini-Med (837*) de Google, ne sont formés qu'à partir de données publiques anonymisées ou pseudonymisées sur les patients – mais des recherches supplémentaires sont nécessaires pour évaluer les risques associés à cette pratique. Des approches de formation préservant la confidentialité ou des données synthétiques peuvent contribuer à résoudre ce problème, comme indiqué dans [la section 3.4.3. Méthodes techniques de protection de la vie privée](#).

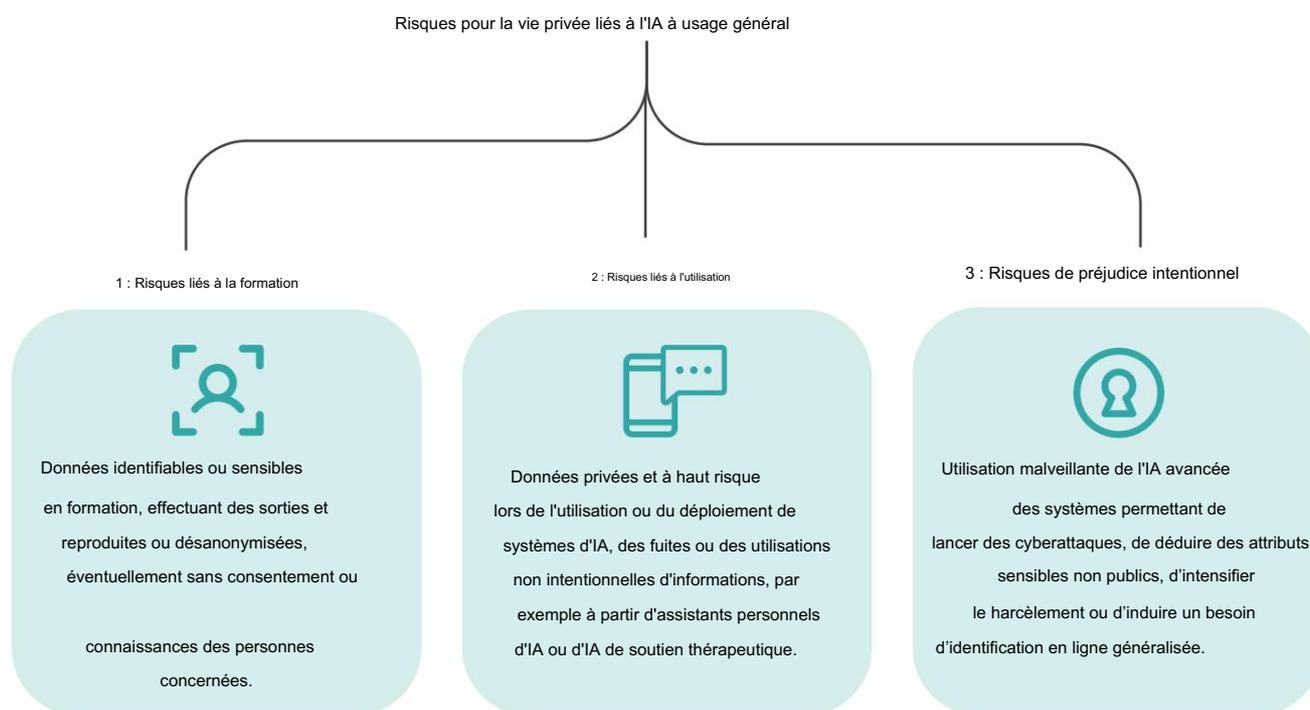


Figure 2.11 : Les risques pour la vie privée liés à l'IA à usage général se répartissent en trois groupes de risques : 1. Risques liés à la formation : risques associés à la formation sur des données sensibles, 2. Risques liés à l'utilisation : risques liés à la manipulation d'informations sensibles lors de l'utilisation de l'IA à usage général, et 3. Risques de préjudice intentionnel : risques liés aux acteurs malveillants appliquant l'IA à usage général pour compromettre la vie privée des individus. Source : International AI Safety Report.

Les informations utilisées lors de l'application de l'IA à usage général peuvent être divulguées, telles que les données privées utilisées pour personnaliser les réponses (« risques d'utilisation »). Les modèles d'IA à usage général n'ont pas connaissance de l'actualité qui se produit après leur formation ni d'informations privées non incluses dans les données de formation. Pour remédier à ce problème, il est courant de fournir des informations de contextualisation pertinentes

aux systèmes d'IA pendant l'utilisation grâce à ce que l'on appelle la « génération augmentée de récupération » (RAG) (838, 839, 840). Ce processus peut également permettre des réponses personnalisées à l'aide de données personnelles privées, par exemple avec des assistants personnels IA sur les téléphones (4*, 841*). Il peut également être utilisé pour inclure des informations externes, telles que des résultats de recherche sur le Web (85*), dans le contexte utilisé pour fournir une réponse. Ceux-ci peuvent être combinés ; par exemple, un outil d'assistance IA pour les soins de santé peut inclure ou accéder à des données médicales sensibles dossiers sur une personne, puis recherchez sur le Web ou dans des bases de données médicales des informations pertinentes

avant de fournir une réponse pour soutenir un clinicien. Si l'utilisation de données privées sur l'appareil peut rendre l'IA à usage général plus utile, elle peut créer des risques supplémentaires de fuite de ces données. Les risques de fuite d'informations vers des tiers augmentent considérablement lorsque les données (ou les informations issues des données) quittent un appareil (842, 843), bien que les approches de cybersécurité puissent minimiser ces risques (844). En pratique, équilibrer la confidentialité, la transparence de l'utilisateur et l'utilité pour le consommateur dans ce contexte est un défi difficile ; des approches techniques pour équilibrer cela existent (voir [3.4.3. Méthodes techniques pour la confidentialité](#)), mais il est également important de trouver des approches politiques qui protègent les droits, permettent la transparence et créent la confiance pour le partage des données afin de promouvoir l'innovation.

Les systèmes d'IA à usage général pourraient permettre une augmentation des atteintes à la vie privée par des acteurs malveillants (« risques de préjudice intentionnel »). Il existe de nombreux scénarios relatifs au risque de confidentialité dans lesquels des utilisateurs malveillants peuvent exploiter les capacités accrues de traitement de l'information de l'IA. Par exemple, des capacités de recherche fine sur Internet, telles que la puissante recherche d'images inversées ou des formes de détection de style d'écriture, permettent d'identifier et de suivre des individus sur des plateformes en ligne, et de déduire des caractéristiques personnelles sensibles (483*, 845) (telles que le sexe, la race, les conditions médicales ou les préférences personnelles), érodant encore davantage la vie privée des individus (846). Les LLM peuvent permettre des recherches plus efficaces et plus efficaces d'informations sensibles dans les données. La détection, la rédaction ou la suppression d'informations personnellement identifiables ne suffisent pas à atténuer complètement l'inférence de contenu personnel sensible : de nombreux attributs d'utilisateur, tels que des préférences sexuelles détaillées ou des habitudes spécifiques de consommation de drogues, peuvent souvent encore être trouvés à partir de données « rédigées » (847), bien que les systèmes d'IA puissent également être utiles pour soutenir la surveillance et la suppression d'informations sensibles en ligne. Ces risques peuvent survenir dans de nombreux contextes et peuvent entraîner un traitement non autorisé à grande échelle des données personnelles. Cela inclut l'usage de la capacité des systèmes d'IA à usage général à déduire des informations privées en fonction des entrées du modèle (316*, 483*). Au-delà de l'analyse et de la recherche, le contenu d'IA à usage général généré à l'aide de données privées, comme les deepfakes non consentis, peut être utilisé pour manipuler ou nuire aux individus. Cela soulève des inquiétudes quant aux dommages causés par l'utilisation malveillante de données personnelles et à l'érosion de la confiance dans le contenu en ligne (voir [2.1.1. Dommages causés aux individus par le biais de faux contenus](#) pour une discussion plus détaillée).

Depuis la publication du rapport intermédiaire, l'importance et les capacités accrues de l'IA à usage général ont conduit à son utilisation accrue dans des contextes sensibles et à un examen ultérieur de ses éventuelles violations des lois sur la vie privée. L'IA à usage général est désormais plus courante dans les contextes où les données sont sensibles, comme les appareils personnels avec assistants intelligents (4*, 841*) et les soins de santé (848*). Jusqu'à présent, aucun grand fournisseur d'IA n'a signalé de fuites d'informations commerciales ou d'utilisateurs très médiatisées, ce qui est significatif étant donné que la divulgation des violations de données personnelles est obligatoire dans la plupart des juridictions. En outre, les chercheurs n'ont pas trouvé de preuve de violations explicites de la vie privée à l'aide de l'IA à usage général. Cependant, contrairement à d'autres préjudices, certaines formes de violations de la vie privée peuvent rester cachées pendant de longues périodes. Par exemple, les atteintes à la vie privée résultant d'une formation sur des données sensibles peuvent ne pas se manifester pendant une période prolongée après la formation, car le temps entre la collecte ou l'utilisation des données et le déploiement ultérieur d'un système d'IA peut être considérable. Les régulateurs appliquent de plus en plus les lois sur la vie privée pour protéger les consommateurs des entreprises qui utilisent l'IA sans contrôles ou garanties de confidentialité (849, 850). Entre-temps, de nouvelles modalités d'interactions avec

L'IA à usage général crée de nouveaux risques pour la vie privée. Par exemple, les modèles de génération de vidéos de haute qualité (851*) peuvent être capables de mémoriser des informations vidéo (comme les visages des élèves dans des salles de classe diffusées en direct) ou d'être utilisés pour exploiter la vie privée en raisonnant sur des données vidéo (852*) ou par l'identification du locuteur (3*) (par exemple, en utilisant l'IA à usage général pour observer les individus et prendre automatiquement des notes sur leur comportement). D'autres préoccupations concernant la vie privée découlant des conséquences en aval de l'IA à usage général ont également émergé. Par exemple, à l'avenir, il pourrait être nécessaire de différencier les humains de l'IA à usage général compétente en ligne, ce qui pourrait rendre l'identification de masse et la surveillance en ligne ultérieure plus probables (853).

Français Les principales lacunes en matière de données probantes concernant la confidentialité portent sur les cas où des informations privées peuvent être divulguées involontairement, sur les moyens de les empêcher et sur les conséquences sociétales que l'IA à usage général pourrait avoir sur la confidentialité. Il est difficile d'évaluer dans quelle mesure l'IA à usage général mémorise ses données d'entraînement et dans quelle mesure elle est susceptible de régurgiter ces données (171, 831). De même, les recherches en cours visent à déterminer dans quelle mesure l'IA à usage général peut ou veut garder privées les informations fournies pendant son utilisation (847). Plus généralement, des recherches sont nécessaires sur les conséquences à long terme pour la confidentialité qui peuvent découler de l'utilisation généralisée de l'IA à usage général, notamment les risques que des acteurs déduisent correctement des informations sensibles sur des individus à l'aide de l'IA à usage général (483*), les risques d'une surveillance de masse renforcée (439, 483*) et les conséquences de l'IA à usage général répandue sur la confidentialité et l'identité (853).

Pour les décideurs politiques travaillant sur la protection de la vie privée, les principaux défis consisteront à évaluer l'ampleur et l'impact des violations de la vie privée par et via l'IA à usage général. Savoir quand et comment la vie privée est violée est intrinsèquement difficile, tant pour les individus que pour les décideurs politiques (854), avec des risques couvrant de multiples aspects du développement et de l'utilisation (résumés dans la figure 2.11). Souvent, des données personnelles sont traitées sans autorisation ou des informations sensibles sont divulguées sans préjudice notable pour l'individu à court terme, ce qui rend difficile d'obtenir un soutien pour traiter les risques de confidentialité de manière préventive (855*). Lorsque des informations sensibles sont divulguées, il peut également être difficile de vérifier où la fuite s'est produite dans les systèmes techniques sous-jacents à l'IA à usage général, car les données sont souvent traitées sur plusieurs appareils ou dans différentes parties de la chaîne d'approvisionnement. Pour les décideurs politiques, ces deux éléments peuvent rendre extrêmement difficile de voir l'ampleur ou la portée des violations de la vie privée, ce qui peut à son tour compliquer la détermination du type et de l'ampleur appropriés de l'intervention. Équilibrer les risques pour la vie privée avec l'utilité des systèmes d'IA à usage général sera difficile mais possible, et des recherches supplémentaires seront nécessaires pour évaluer les risques et minimiser les dommages.

Pour les pratiques de gestion des risques liées à la confidentialité, voir :

- [3.4.2. Suivi et intervention](#)
- [3.4.3. Méthodes techniques de protection de la vie privée](#)

2.3.6. Risques de violation du droit d'auteur

INFORMATIONS CLÉS

- L'utilisation de grandes quantités de données pour la formation de modèles d'IA à usage général a entraîné des préoccupations liées aux droits sur les données et à la propriété intellectuelle. La collecte de données et la génération de contenu peuvent impliquer une variété de lois sur les droits des données, qui varient selon les juridictions et peuvent faire l'objet de litiges actifs. Compte tenu de l'incertitude juridique entourant les pratiques de collecte de données, les entreprises d'IA partagent moins d'informations sur les données qu'elles utilisent. Cette opacité rend plus difficile la recherche sur la sécurité de l'IA par des tiers.
- La création de contenu par l'IA remet en question les systèmes traditionnels de consentement, de rémunération et de contrôle des données. Les lois sur la propriété intellectuelle sont conçues pour protéger et promouvoir l'expression créative et l'innovation. L'IA à usage général apprend des œuvres d'expression créative et peut les créer.
- Les chercheurs développent des outils et des méthodes pour atténuer les risques potentiels. Les lois sur la violation des droits d'auteur et d'autres droits relatifs aux données restent toutefois peu fiables. Il existe également des outils limités pour rechercher et filtrer les données de formation à grande échelle en fonction de leurs licences, du consentement explicite des créateurs ou d'autres critères juridiques et éthiques.
- Depuis le rapport intermédiaire (mai 2024), les détenteurs de droits sur les données ont rapidement restreint l'accès à leurs données. Cela empêche les développeurs d'IA d'utiliser ces données pour former leurs modèles, mais entrave également l'accès aux données à des fins de recherche, de bien social ou à des fins non liées à l'IA.
- Les décideurs politiques doivent relever le défi de permettre un accès responsable et conforme aux lois aux données sans décourager le partage des données et l'innovation. Les outils techniques permettant d'évaluer, de tracer, de filtrer et d'attribuer automatiquement des licences aux données pourraient faciliter cette tâche, mais les outils actuels ne sont pas suffisamment évolutifs et efficaces.

Définitions clés

- **Propriété intellectuelle** : Créations de l'esprit sur lesquelles des droits légaux peuvent être accordés, y compris œuvres littéraires et artistiques, symboles, noms et images.
- **Droit d'auteur** : Une forme de protection juridique accordée aux créateurs d'œuvres originales, leur donnant des droits exclusifs d'utilisation, de reproduction et de distribution de leurs œuvres.
- **Marque déposée** : Un symbole, un mot ou une phrase légalement enregistré ou établi par l'usage pour représenter une entreprise ou produit, le distinguant des autres sur le marché.
- **Droits d'image** : droits qui protègent l'image, la voix, le nom ou d'autres aspects identifiables d'une personne contre toute utilisation commerciale non autorisée.
- **Fair use** : doctrine juridique américaine qui fournit une défense contre la violation du droit d'auteur revendications d'utilisation limitée de matériel protégé par le droit d'auteur sans autorisation à des fins telles que la critique, le commentaire, le reportage d'actualité, l'éducation et la recherche. Certains autres pays autorisent des droits d'utilisation similaires sous le nom de « fair handling ».

- Exploration Web : Utilisation d'un programme automatisé, souvent appelé robot d'exploration ou bot, pour naviguer sur le Web. web, aux fins de collecte de données à partir de sites Web.

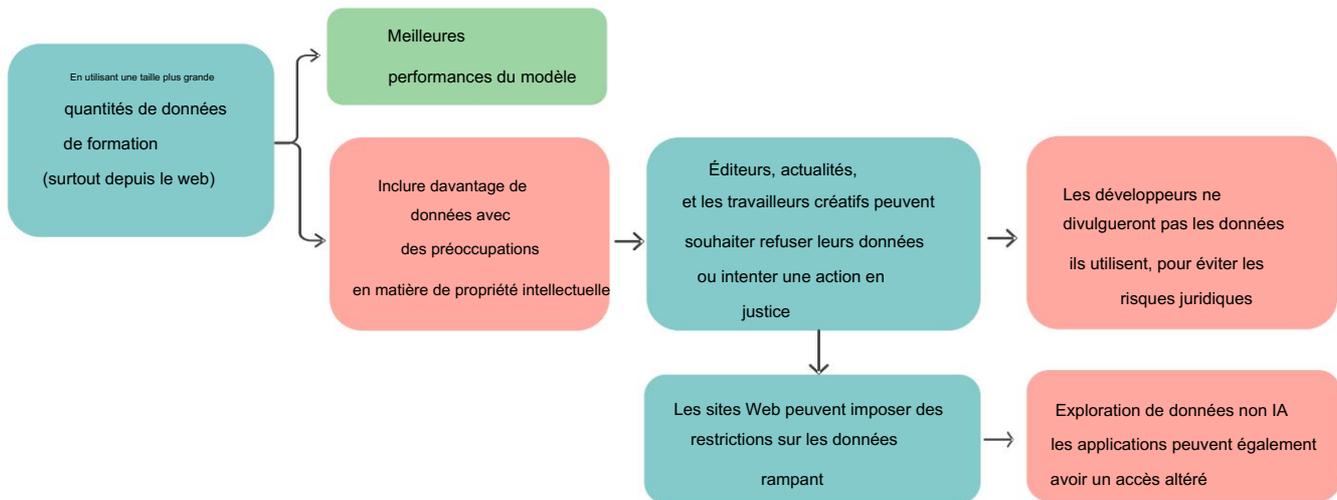


Figure 2.12 : Les avantages liés à l'utilisation de grandes quantités de données d'apprentissage peuvent avoir des conséquences en cascade sur la transparence des données, l'exploration du Web et les normes de partage d'informations sur le Web. Source : International AI Safety Report.

L'IA à usage général s'entraîne sur de grandes collections de données, ce qui peut impliquer une variété de lois sur les droits des données, notamment la propriété intellectuelle, la confidentialité, les marques déposées et les droits à l'image/à la ressemblance.

L'IA à usage général est formée sur de grandes collections de données, souvent provenant en partie d'Internet.

Ils peuvent être utilisés pour générer du texte, des images, du son ou des vidéos qui peuvent parfois imiter le contenu sur lequel ils ont été formés. Dans le cas de la collecte de données (entrées) et de la génération de données (sorties), ces systèmes peuvent impliquer divers droits et lois sur les données (voir la figure 2.12). Par exemple, si les données de formation de l'IA contiennent des informations personnellement identifiables, cela peut engendrer des problèmes de confidentialité.

De même, les ensembles de données de formation provenant du Web contiennent souvent des éléments protégés par des droits d'auteur, ce qui implique des lois sur le droit d'auteur et la propriété intellectuelle (836, 856). Si des marques sont capturées dans les données, des marques commerciales peuvent également être impliquées. Dans certaines juridictions, des personnes célèbres figurant dans les données de formation peuvent avoir des droits d'image (857). Les lois régissant ces droits sur les données peuvent également varier selon les juridictions et, en particulier dans le cas de l'IA, certaines font l'objet de litiges actifs.

Les lois sur le droit d'auteur visent à protéger l'expression créative ; l'IA à usage général apprend et génère du contenu ressemblant à l'expression créative. Les lois sur le droit d'auteur visent à protéger et à encourager l'expression écrite et créative (858, 859), principalement sous la forme d'œuvres littéraires (y compris les logiciels), d'arts visuels, de musique, d'enregistrements sonores et d'œuvres audiovisuelles. Elles accordent aux créateurs d'œuvres originales le droit exclusif de copier, distribuer, adapter et exécuter leur propre travail. L'utilisation non autorisée par des tiers de données protégées par le droit d'auteur est autorisée dans certaines juridictions et circonstances : par exemple sur la base de l'exception d'« usage équitable » aux États-Unis (860), de l'exception d'« exploration de textes et de données » dans l'UE (861), de la loi modifiée sur le droit d'auteur au Japon (862), de la loi israélienne sur le droit d'auteur (863) et de la loi sur le droit d'auteur de 2021 à Singapour (864). Dans chaque juridiction, il existe différentes lois relatives à (a) la licéité des pratiques de collecte de données (par exemple, le scraping de données), (b) l'utilisation des données (par exemple pour la formation d'IA, de systèmes commerciaux ou non commerciaux) et (c) si

Les résultats de modèles qui semblent similaires à du matériel protégé par le droit d'auteur constituent une infraction. Aux États-Unis, ces questions font l'objet de litiges actifs (865, 866, 867, 868, 869), par exemple dans des affaires telles que le New York Times contre OpenAI et Microsoft. De nombreux problèmes liés à la création et à l'utilisation d'ensembles de données tout au long du cycle de vie de l'ensemble de données rendent les préoccupations en matière de droits d'auteur pour la formation de modèles d'IA très compliquées (870). Les questions pertinentes incluent notamment si les ensembles de données ont été assemblés spécifiquement pour l'apprentissage automatique ou à l'origine à d'autres fins (871), si la violation potentielle s'applique aux entrées et/ou aux sorties de modèles (872, 873, 874) et à quelle juridiction l'affaire relève, entre autres (481). Des questions se posent également quant à savoir qui est responsable de la violation ou des sorties de modèles préjudiciables (développeurs, utilisateurs ou autres acteurs) (875). Bien que les développeurs puissent utiliser des stratégies techniques pour atténuer les risques de violation du droit d'auteur à partir des sorties de modèles, ces risques sont difficiles à éliminer entièrement (876, 877).

Les systèmes d'IA à usage général peuvent avoir un impact sur les économies créatives et d'édition. À mesure que les systèmes d'IA à usage général deviennent plus performants, ils ont de plus en plus le potentiel de perturber les marchés du travail, et en particulier les industries créatives (662, 707) (voir également [2.3.1. Risques liés au marché du travail](#)). Les décisions judiciaires en attente concernant la violation du droit d'auteur dans la phase de formation de l'IA peuvent affecter la capacité des développeurs d'IA à usage général à créer des modèles puissants et performants en limitant leur accès aux données de formation (836, 856, 878). Elles peuvent également avoir un impact sur la capacité des créateurs de données à limiter l'utilisation de leurs données, ce qui peut décourager l'expression créative. Par exemple, les éditeurs de presse et les artistes ont exprimé leur inquiétude quant au fait que leurs clients pourraient utiliser des systèmes d'IA générative pour leur fournir un contenu similaire. Dans les domaines de l'information, de l'art et du divertissement, l'IA générative peut souvent produire des versions paraphrasées, abstraites ou résumées du contenu sur lequel elle a été formée. Si les utilisateurs accèdent aux informations via des résumés d'IA générative plutôt qu'à partir de sites de médias, cela pourrait réduire les revenus d'abonnement et de publicité pour les éditeurs d'origine. Les abonnements réduits peuvent être assimilés à des dommages-intérêts pour atteinte aux droits d'auteur.

L'incertitude juridique entourant les pratiques de collecte de données a découragé la transparence sur les données collectées ou utilisées par les développeurs d'IA à usage général, ce qui rend plus difficile la recherche sur la sécurité de l'IA par des tiers. Les chercheurs indépendants en IA peuvent plus facilement comprendre les risques et les dangers potentiels d'un système d'IA à usage général s'il existe une transparence sur les données sur lesquelles il a été formé (879). Par exemple, il est beaucoup plus facile de quantifier le risque qu'un modèle génère des informations biaisées, protégées par des droits d'auteur ou privées si le chercheur sait sur quelles sources de données il a été formé.

Cependant, ce type de transparence fait souvent défaut aux principaux développeurs d'IA à usage général (880). La crainte du risque juridique, notamment en cas de violation du droit d'auteur, dissuade les développeurs d'IA de divulguer leurs données de formation (881).

L'infrastructure permettant de rechercher et de filtrer les données légalement autorisées est sous-développée, ce qui rend difficile pour les développeurs de se conformer à la loi sur le droit d'auteur. La possibilité d'utiliser des œuvres protégées par le droit d'auteur dans le cadre de données de formation sans autorisation appropriée est un domaine de litige actif. Les outils permettant de rechercher et d'identifier les données disponibles sans problème de droit d'auteur sont limités. Par exemple, des travaux récents montrent qu'environ 60 % des ensembles de données populaires dans les référentiels de données en accès libre les plus utilisés contiennent des informations de licence incorrectes ou manquantes (481). De même, les outils actuels permettant de distinguer les données libres de droits dans les scrapings Web présentent des limites (856, 878). Cependant,

les praticiens développent de nouvelles normes pour la documentation des données et de nouveaux protocoles pour que les créateurs de données signalent leur consentement à l'utilisation dans la formation des modèles d'IA (882, 883).

Depuis la publication du rapport intermédiaire, les luttes juridiques et techniques autour des données se sont intensifiées et les recherches suggèrent qu'il reste difficile d'empêcher complètement les modèles de générer du contenu protégé par le droit d'auteur en utilisant des mesures d'atténuation techniques. De nombreuses organisations, y compris les développeurs d'IA, utilisent des robots automatiques appelés « robots d'exploration Web » qui naviguent sur le Web et copient le contenu.

Les sites Web souhaitent souvent que leur contenu soit lu par des robots d'exploration qui dirigeront le trafic humain vers eux (tels que les robots d'exploration des moteurs de recherche) mais laissé tranquille par des robots d'exploration qui copieront leurs données pour entraîner des outils concurrents (par exemple des modèles d'IA qui déplaceront leur trafic). Les sites Web peuvent indiquer leurs préférences aux robots d'exploration dans leur code, y compris si et par qui ils souhaitent être explorés. Ils peuvent également utiliser des technologies qui tentent d'identifier et de bloquer les robots d'exploration. Depuis mai 2024, des preuves sont apparues selon lesquelles les sites Web ont érigé davantage de barrières aux robots d'exploration de la part des développeurs d'IA (836, 884, 885). Ces mesures sont déclenchées par l'incertitude quant à savoir si les robots d'exploration des développeurs d'IA respecteront les signaux de préférence des sites Web. À la recherche de solutions, le Bureau européen de l'IA élabore un code de pratique de reporting de transparence pour les développeurs d'IA à usage général (886), et le National Institute of Standards and Technology (NIST) des États-Unis a publié un cadre de gestion des risques liés à l'IA (887). En outre, un nombre croissant de travaux étudient la capacité des chercheurs à extraire des informations d'un modèle entraîné ou à détecter sur quoi un modèle a été entraîné. Cependant, ces méthodes, appliquées aux modèles d'IA à usage général, présentent encore des défis fondamentaux qui pourraient ne pas être facilement surmontés dans un avenir proche (831, 832, 888, 889, 890, 891, 892).

L'augmentation des obstacles à l'accès au contenu Web peut entraver la collecte de données, y compris pour les applications non liées à l'IA. L'augmentation des restrictions sur l'exploration du Web fait que les données de la plus haute qualité et bien entretenues sont moins disponibles, en particulier pour les organisations moins bien financées (836, 856, 878). La diminution de la disponibilité des données peut avoir des répercussions sur la concurrence, sur la diversité et la factualité des données de formation, ainsi que sur la capacité des régions mal desservies à développer leurs propres applications d'IA compétitives. Si les grandes entreprises d'IA peuvent se permettre d'acheter des licences de données ou simplement développer des robots d'exploration plus puissants pour accéder aux données restreintes, l'augmentation des restrictions aura des externalités négatives pour les autres utilisations (y compris de nombreuses utilisations bénéfiques) des robots d'exploration du Web. De nombreux secteurs dépendent des robots d'exploration : recherche sur le Web, catalogues de produits/prix, études de marché, publicité, archives Web, recherche universitaire, outils d'accessibilité et même applications de sécurité. L'accès de ces secteurs aux données est de plus en plus limité en raison des obstacles érigés pour empêcher les grands développeurs d'IA d'utiliser les données pour la formation. Enfin, ces défis liés aux robots d'exploration peuvent persister, même lorsque les litiges relatifs aux droits d'auteur sont résolus.

Les outils permettant de découvrir, d'étudier et d'attribuer automatiquement des licences aux données font défaut. Des outils standardisés sont nécessaires aux créateurs et aux utilisateurs de données pour évaluer les restrictions ou les limites d'un ensemble de données, pour estimer la valeur des données, pour les attribuer automatiquement sous licence à grande échelle et pour suivre leur utilisation en aval (465, 481, 856, 878). Sans ces outils, le marché s'est jusqu'à présent appuyé sur des contrats ad hoc et personnalisés, sans processus de licence clair pour les petits créateurs de données. Associées au manque de transparence des données de la part des développeurs individuels, ces lacunes entravent le développement d'un système de gestion efficace et

marché des données structurées. En substance, le Web est une source de données relativement désordonnée et non structurée.

Sans de meilleurs outils pour l'organiser, les développeurs auront du mal à éviter de se former sur des données qui peuvent engendrer des problèmes juridiques ou éthiques.

Les méthodes qui atténuent le risque de violation du droit d'auteur dans les modèles sont sous-développées et nécessitent davantage de recherche. Les grands modèles peuvent mémoriser ou rappeler certaines des données sur lesquelles ils ont été formés, ce qui leur permet de les reproduire lorsqu'ils y sont invités. Par exemple, des sections des livres Harry Potter sont mémorisées dans des modèles de langage courant (893*). Cela est souhaitable dans certains cas (par exemple, rappeler des faits), mais indésirable dans d'autres, car cela peut conduire les modèles à générer et à redistribuer du matériel protégé par le droit d'auteur, des informations privées ou du contenu sensible trouvé sur le Web. Il existe de nombreuses approches pour atténuer ce risque (voir également 3.4.3. [Méthodes techniques de confidentialité](#)). L'une consiste à détecter si un modèle a été formé sur ou a mémorisé certains contenus indésirables, ce qui lui permettrait également de les régénérer (831, 832, 888, 889). C'est ce qu'on appelle la « recherche de mémorisation » ou la « recherche d'inférence d'appartenance ». Les chercheurs peuvent également examiner si les résultats du modèle peuvent être attribués directement à certains points de données d'entraînement (877, 890). Une autre méthode consiste à utiliser des filtres qui détectent quand un modèle génère du contenu substantiellement similaire à du matériel protégé par le droit d'auteur. Cependant, il reste difficile, sur le plan conceptuel et technique, de tester si les générations sont substantiellement similaires au contenu protégé par le droit d'auteur sur lequel le modèle a été formé (891, 894). Enfin, les chercheurs explorent des méthodes pour supprimer les informations que les modèles ont déjà apprises, appelées « désapprentissage automatique » (821, 895, 896, 897, 898). Cependant, il se peut que cette solution ne soit pas viable, robuste ou pratique à long terme (892, 897, 898). Par exemple, le désapprentissage automatique ne parvient souvent pas à supprimer complètement les informations ciblées d'un modèle et, ce faisant, il peut déformer les autres capacités du modèle de manière imprévue, ce qui le rend peu attrayant pour les développeurs d'IA commerciale (892, 895, 897, 898).

Les décideurs politiques sont confrontés au défi de permettre la protection de la propriété intellectuelle et d'autres droits sur les données, tout en créant un environnement qui encourage le partage des données pour promouvoir l'innovation. Ces défis sont exacerbés par les nombreuses lois applicables, qui varient selon les juridictions ou font l'objet de litiges actifs. Ils sont également compliqués par le manque de transparence des données existantes dans le développement de l'IA et par la vague de détenteurs de droits sur les données qui recourent à leurs propres mesures pour protéger leurs données. Dans l'ensemble, l'écosystème du Web et la chaîne d'approvisionnement des données évoluent rapidement en réponse à l'IA, avec ou sans interventions juridiques. Ces tendances démontrent les défis à relever pour encourager une plus grande transparence et développer des solutions techniques qui permettent un marché plus sain des données. Sans de telles solutions, un manque de

La transparence dans l'utilisation des données entravera la recherche sur la sécurité de l'IA, aura un impact négatif sur les économies créatives et encouragera davantage de protectionnisme des données, avec des conséquences au-delà du développement de l'IA.

Pour les pratiques de gestion des risques liées au droit d'auteur, voir :

- [3.3. Identification et évaluation des risques](#)
- [3.4.1. Former des modèles plus fiables](#)
- [3.4.3. Méthodes techniques de protection de la vie privée](#)

2.4. Impact des modèles d'IA polyvalents à pondération ouverte sur les risques liés à l'IA

INFORMATIONS CLÉS

- La manière dont un modèle d'IA est publié est un facteur important dans l'évaluation des risques qu'il présente. Il existe toute une gamme d'options de publication de modèles, de la version totalement fermée à la version totalement ouverte, qui impliquent toutes des compromis entre risques et avantages. Les modèles à pondération ouverte (ceux dont les pondérations sont rendues publiques pour téléchargement) représentent un point clé de ce spectre.
- Les modèles de pondération ouverts facilitent la recherche et l'innovation, mais permettent également des utilisations malveillantes et la perpétuation de certaines failles. Les pondérations ouvertes permettent aux communautés de recherche mondiales de faire progresser les capacités et de remédier aux failles des modèles en leur fournissant un accès direct à un composant d'IA essentiel dont le développement indépendant est extrêmement coûteux pour la plupart des acteurs. Cependant, la publication ouverte des pondérations des modèles pourrait également présenter des risques de facilitation d'utilisation malveillante ou erronée ou de perpétuation des failles et des biais des modèles.
- Une fois que les pondérations des modèles sont disponibles pour téléchargement public, il n'y a aucun moyen de mettre en œuvre une retrait en masse de toutes les copies existantes du modèle. En effet, divers acteurs auront fait leurs propres copies. Même si elles sont retirées des plateformes d'hébergement, les versions téléchargées existantes sont faciles à distribuer hors ligne. Par exemple, des modèles de pointe tels que Llama-3.1-405B peuvent tenir sur une clé USB.
- Depuis le rapport intermédiaire (mai 2024), un consensus de haut niveau s'est dégagé sur le fait que les risques posés par une plus grande ouverture de l'IA devraient être évalués en termes de « risque marginal ». Il s'agit du risque supplémentaire associé à la diffusion ouverte de l'IA, par rapport aux risques posés par les modèles fermés ou la technologie existante.
- Qu'un modèle soit ouvert ou fermé, des approches d'atténuation des risques doivent être mises en œuvre tout au long du cycle de vie de l'IA, y compris lors de la collecte de données, de la préformation du modèle, du réglage fin et des mesures post-publication. L'utilisation de plusieurs mesures d'atténuation peut renforcer les interventions imparfaites.
- Un des principaux défis pour les décideurs politiques est de comprendre les lacunes en matière de données probantes concernant les impacts positifs et négatifs potentiels de la publication ouverte des poids sur la concentration du marché et la concurrence. Les effets varieront probablement en fonction du degré d'ouverture du modèle (par exemple, si la publication est sous une licence open source), du niveau de marché discuté (c'est-à-dire la concurrence entre les développeurs d'IA à usage général et les développeurs d'applications en aval) et de l'ampleur de l'écart entre les concurrents. • Un autre défi majeur pour les décideurs politiques est de reconnaître les limites techniques de la publication ouverte des poids sur la concentration du marché et la concurrence.

certaines interventions politiques pour les modèles ouverts. Par exemple, des exigences telles que le tatouage numérique robuste pour les modèles d'IA génératifs à pondération ouverte sont actuellement irréalisables, car il existe des limitations techniques à la mise en œuvre des tatouages numériques qui ne peuvent pas être supprimées.

Définitions clés

- Interface de programmation d'application (API) : un ensemble de règles et de protocoles qui permet l'intégration et la communication entre les systèmes d'IA et d'autres applications logicielles.
- Risque marginal : le risque supplémentaire introduit par un modèle ou un système d'IA à usage général par rapport à une référence pertinente, comme un risque comparable posé par une technologie non IA existante.
- Modèle à pondération ouverte : un modèle d'IA dont les pondérations sont disponibles publiquement en téléchargement, comme Llama ou Stable Diffusion. Les modèles à pondération ouverte peuvent être, mais ne sont pas nécessairement, open source.
- Modèle open source : modèle d'IA mis à disposition du public en téléchargement sous une licence open source. La licence open source accorde la liberté d'utiliser, d'étudier, de modifier et de partager le modèle à toutes fins utiles. Il subsiste un certain désaccord quant aux composants du modèle (poids, code, données d'entraînement) et à la documentation qui doivent être accessibles au public pour que le modèle soit qualifié d'open source.
- Pondérations : paramètres de modèle qui représentent la force de connexion entre les nœuds d'un réseau neuronal. Les pondérations jouent un rôle important dans la détermination de la sortie d'un modèle en réponse à une entrée donnée et sont mises à jour de manière itérative pendant l'apprentissage du modèle pour améliorer ses performances.
- Deepfake : un type de faux contenu généré par l'IA, composé de contenu audio ou visuel, qui déforme les faits et présente des personnes réelles comme faisant ou disant quelque chose qu'elles n'ont pas réellement fait ou dit.

Cette section se concentre principalement sur les avantages et les risques des modèles d'IA à usage général avec des pondérations de modèle largement disponibles. Les pondérations de modèle, également appelées paramètres, sont les nombres utilisés pour spécifier comment l'entrée (par exemple, un texte décrivant une image) est transformée en sortie (par exemple, l'image elle-même).

Ces pondérations sont mises à jour de manière itérative au cours de l'entraînement du modèle afin d'améliorer les performances du modèle sur les tâches pour lesquelles il est entraîné (voir [1.1. Comment l'IA à usage général est-elle développée](#)). Si la réalisation de tous les avantages de l'ouverture de l'IA nécessite une ouverture plus poussée que le simple partage des pondérations du modèle (par exemple, cela nécessite le partage des données d'entraînement, du code d'entraînement, de la documentation, etc.), de nombreux risques associés à la diffusion ouverte des modèles découlent du fait que les pondérations du modèle sont rendues accessibles au public (899). En conséquence, les modèles à pondération ouverte font l'objet de nombreux travaux politiques.

La différence entre les modèles « open-weight » et les modèles « open source » peut être déroutante.

« Open-weight » signifie que les pondérations du modèle sont disponibles pour téléchargement public, comme avec Llama, Mixtral ou Hunyuan-Large. Les modèles à pondération ouverte peuvent être, mais ne sont pas nécessairement, open source. La classification « open source » exige que l'accès au modèle soit protégé par une licence open source qui accorde la liberté légale à quiconque d'utiliser, d'étudier, de modifier et de partager le modèle à n'importe quelle fin. Les licences open source sont importantes pour réaliser les avantages de l'ouverture de l'IA : elles favorisent l'innovation et contrent la concentration du marché par les grandes entreprises technologiques en permettant aux développeurs en aval d'utiliser, d'étudier et de modifier des modèles ouverts sans avoir à demander la permission, et d'intégrer ces modèles dans des produits qu'ils peuvent mettre sur le marché. Cela inclut des avantages pour les acteurs à faibles ressources qui ne pourraient autrement pas avoir accès aux pondérations des modèles, car leur formation à partir de zéro est coûteuse. Bien que la licence open source soit essentielle pour ouvrir

classification du modèle source, il subsiste un certain désaccord quant à la mesure dans laquelle les différents composants (poids, code, données de formation) et la documentation doivent être accessibles au public pour que le modèle soit qualifié d'open source.

Il existe également un éventail d'options de diffusion de modèles, allant de la fermeture totale à l'ouverture totale, qui impliquent toutes des compromis entre les risques et les avantages (voir tableau 2.5).

- Les modèles entièrement ouverts sont des modèles open source pour lesquels les pondérations, le code complet, les données de formation et les autres documents (par exemple sur le processus de formation du modèle) sont mis à disposition du public, sans restrictions de modification, d'utilisation et de partage. En général, la publication entièrement ouverte du modèle facilite la recherche et l'innovation à plus grande échelle, mais augmente les risques d'utilisation malveillante en permettant aux acteurs malveillants de contourner facilement les restrictions de sécurité et de modifier le modèle à des fins nuisibles, et en augmentant la probabilité que les défauts du modèle prolifèrent en aval dans les versions et applications modifiées du modèle si les utilisateurs en aval ne mettent pas à jour de manière proactive la version du modèle qu'ils utilisent. • Les poids et le code des modèles entièrement fermés sont exclusifs, pour un usage interne uniquement. Cela signifie que

Les acteurs externes ne peuvent pas utiliser le modèle à mauvais escient et les failles sont moins susceptibles de proliférer en aval et peuvent être corrigées une fois découvertes. Cependant, avec des modèles fermés, il est également plus difficile pour les développeurs externes de découvrir les risques d'utilisation abusive, les failles et d'utiliser le modèle à des fins plus larges. innovation et recherche.

- Les modèles partiellement ouverts partagent une certaine combinaison de poids, de code et de données sous diverses conditions. Les licences ou les contrôles d'accès, dans le but de trouver un équilibre entre les avantages de l'ouverture et l'atténuation des risques et les préoccupations liées à la propriété intellectuelle. Par exemple, OpenAI fournit un accès public à son modèle GPT-4o via une interface appelée ChatGPT qui permet aux utilisateurs d'interroger le système et de récupérer des réponses sans accéder au modèle lui-même. Ce type d'« accès partiel aux requêtes » permet au public d'utiliser le modèle et d'étudier son comportement et ses défauts de performance sans fournir un accès direct aux pondérations et au code du modèle. Le coût de cet accès partiel est que les chercheurs en IA externes (universitaires et évaluateurs tiers) n'ont pas accès pour effectuer une analyse plus approfondie de la sécurité du système, et les développeurs en aval ne peuvent pas intégrer librement le modèle dans de nouvelles applications et de nouveaux produits. Certaines licences telles que RAIL (Responsible AI License) énoncent des restrictions contre les utilisations nuisibles du modèle. Les restrictions de licence ne sont que des expressions juridiques et ne constituent aucun obstacle physique à une mauvaise utilisation si le modèle lui-même est disponible en téléchargement public. Certains acteurs peuvent être dissuadés de faire un usage abusif par le risque de responsabilité légale, tandis que d'autres acteurs malveillants peuvent simplement ignorer la condition de licence.

Niveau d'accès	Ce que cela signifie	Exemples	Logiciels traditionnels Analogie
Entièrement fermé	Les utilisateurs ne peuvent pas du tout interagir directement avec le modèle	Flamant (Google)	Algorithmes de trading utilisés par les fonds spéculatifs privés
Accès hébergé	Les utilisateurs ne peuvent interagir que via une application ou une interface spécifique	À mi-parcours (À mi-parcours)	Logiciel grand public Cloud (par exemple Google Docs)
Accès API au modèle	Les utilisateurs peuvent envoyer des requêtes au modèle par programmation, ce qui permet une utilisation dans des applications externes	Claude 3.5 Sonnet (Anthropic)	API basée sur le cloud (par exemple, créateurs de sites Web tels que Squarespace)
Accès API à Réglage fin	Les utilisateurs peuvent affiner le modèle pour leurs besoins spécifiques	GPT-4o (OpenAI)	Logiciel d'entreprise avec API de personnalisation (par exemple Plateforme de développement Salesforce)
Poids ouvert : Poids disponibles À télécharger	Les utilisateurs peuvent télécharger et exécuter le modèle localement	Lama 3 (Méta), Mistral	Logiciel de bureau propriétaire (par exemple Microsoft Word)
Poids, données et Code disponible pour Télécharger avec utilisation Restrictions	Les utilisateurs peuvent télécharger et exécuter le modèle ainsi que le code d'inférence et de formation, mais avoir une certaine licence restrictions sur leur utilisation	FLORAISON (Grande Science)	Logiciel disponible à la source (par exemple Unreal Engine)
Entièrement ouvert : Poids, Données et code Disponible pour Télécharger sans utiliser Restrictions	Les utilisateurs ont une totale liberté de télécharger, d'utiliser et de modifier le modèle, le code complet et les données	GPT-NéoX (EleutherAI)	Logiciels open source (par exemple Mozilla Firefox et Linux)

Tableau 2.5 : Il existe un éventail d'options de partage de modèles allant des modèles entièrement fermés (les modèles sont privés et réservés à un usage exclusif) aux modèles entièrement ouverts et open source (les pondérations, les données et le code des modèles sont librement et publiquement disponibles sans restriction d'utilisation, de modification et de partage). Cette section se concentre sur les trois colonnes les plus à droite. Source : adapté de Bommasani et al., 2024 (900).

Une plus grande ouverture de l'IA présente des avantages, notamment en facilitant l'innovation, en améliorant la sécurité et la surveillance de l'IA, en augmentant l'accessibilité et en permettant aux outils d'IA d'être adaptés à des besoins divers. La formation d'un modèle d'IA à usage général (le processus de production de pondérations de modèle) est extrêmement coûteuse. Par exemple, on estime que la formation du modèle Gemini de Google a coûté 191 millions de dollars en coûts de calcul uniquement (731). Le coût de la formation des calculs pour le modèle d'IA à usage général le plus cher devrait dépasser 1 milliard de dollars d'ici 2027 (27). Les coûts de formation constituent donc un obstacle insurmontable pour de nombreux acteurs (entreprises, universitaires et États) qui souhaitent participer au marché de l'IA à usage général et bénéficier des applications de l'IA. La publication ouverte des pondérations rend l'IA à usage général plus accessible aux acteurs qui, autrement, pourraient ne pas disposer des ressources nécessaires pour les développer de manière indépendante. Cela réduit la dépendance à l'égard des systèmes propriétaires contrôlés par quelques grandes entreprises technologiques (ou potentiellement des États-nations) et permet aux développeurs d'affiner les systèmes existants.

Des pondérations d'IA à usage général pour répondre à des besoins plus diversifiés. Par exemple, les développeurs issus de groupes linguistiques minoritaires peuvent affiner les modèles à pondération ouverte avec des ensembles de données linguistiques spécifiques pour améliorer les performances du modèle dans cette langue. Les modèles peuvent également être affinés plus librement pour mieux fonctionner dans des tâches spécifiques, telles que la rédaction de textes juridiques professionnels, de notes médicales ou d'écriture créative. En outre, une plus grande ouverture permet à une communauté plus large et plus diversifiée de développeurs et de chercheurs d'évaluer les modèles et d'identifier et de remédier aux vulnérabilités, ce qui peut contribuer à une plus grande sécurité de l'IA et à accélérer l'innovation bénéfique en la matière. En général, plus un modèle est ouvert – y compris s'il donne accès à des composants d'IA supplémentaires au-delà des pondérations du modèle, tels que les données d'entraînement, le code, la documentation et l'infrastructure de calcul requise pour utiliser ces modèles – plus les avantages en termes d'innovation et de surveillance de la sécurité sont importants.

Les risques posés par les modèles à pondération ouverte sont en grande partie liés à la possibilité d'une utilisation malveillante ou malavisée (899, 901, 902). Les modèles d'IA à usage général sont à double usage, ce qui signifie qu'ils peuvent être utilisés à bon escient ou à des fins néfastes. Les pondérations des modèles ouverts peuvent potentiellement exacerber les risques d'utilisation abusive en permettant à un large éventail d'acteurs qui ne disposent pas des ressources et des connaissances nécessaires pour créer eux-mêmes un modèle d'exploiter et d'augmenter les capacités existantes à des fins malveillantes et sans surveillance. Bien que les modèles à pondération ouverte et les modèles fermés puissent tous deux comporter des mesures de protection pour refuser les demandes des utilisateurs, ces mesures de protection sont plus faciles à supprimer pour les modèles ouverts. Par exemple, même si un modèle à pondération ouverte comporte des mesures de protection intégrées, telles que des filtres de contenu ou des ensembles de données d'entraînement limités, l'accès aux pondérations du modèle et au code d'inférence permet aux acteurs malveillants de contourner ces mesures de protection (903). En outre, les vulnérabilités des modèles trouvées dans les modèles ouverts peuvent également exposer des vulnérabilités dans les modèles fermés (904*). Enfin, grâce à l'accès aux pondérations des modèles, les acteurs malveillants peuvent également affiner un modèle afin d'optimiser ses performances pour des applications nuisibles (905, 906, 907). Les utilisations malveillantes potentielles comprennent les applications scientifiques à double usage nuisibles, par exemple l'utilisation de l'IA pour découvrir de nouvelles armes chimiques ([2.1.4. Attaques biologiques et chimiques](#)), les cyberattaques ([2.1.3. Cyberinfraction](#)) et la production de faux contenus nuisibles tels que des contenus d'abus sexuels « deepfake » ([2.1.1. Dommages causés aux individus par le biais de faux contenus](#)) et de fausses nouvelles politiques ([2.1.2. Manipulation de l'opinion publique](#)). Comme indiqué ci-dessous, la diffusion d'un modèle à pondération ouverte avec un potentiel d'utilisation malveillante n'est généralement pas réversible, même lorsque ses risques sont découverts ultérieurement.

Il existe également un risque de perpétuer les failles par le biais de versions ouvertes, bien que l'ouverture permette également à beaucoup plus d'acteurs d'effectuer une analyse technique plus approfondie pour repérer ces failles et biais. Lorsque des modèles d'IA à usage général sont publiés ouvertement et intégrés dans une multitude de systèmes et d'applications en aval, tous les défauts de modèle non résolus - tels que les biais et la discrimination intégrés ([2.2.2. Biais](#)), les vulnérabilités aux attaques adverses (904*) ou la capacité de tromper les [systèmes de surveillance](#) post-déploiement en ayant appris à « battre le test » ([2.2.3. Perte de contrôle](#)) - sont également distribués (902). Le même défi se pose pour les modèles fermés, hébergés ou d'accès API, mais pour [ces modèles non téléchargeables](#), l'hôte du modèle peut déployer universellement de nouvelles versions de modèle pour corriger les vulnérabilités et les défauts. Pour les modèles ouverts, les développeurs peuvent mettre à disposition des versions mises à jour, mais rien ne garantit que les développeurs en aval adopteront les mises à jour.

D'autre part, les modèles à pondération ouverte peuvent être examinés et testés plus en profondeur par un plus grand nombre de chercheurs et de développeurs en aval, ce qui permet d'identifier et de corriger davantage de défauts dans les versions futures (908).

Français Depuis la publication du rapport intermédiaire, un consensus de haut niveau s'est dégagé sur le fait que les risques posés par une plus grande ouverture de l'IA devraient être évalués en termes de risque marginal (901, 909, 910). Le « risque marginal » fait référence au risque supplémentaire que représente la diffusion ouverte de l'IA par rapport aux risques posés par les alternatives existantes, telles que les modèles fermés ou d'autres technologies (911). Les études qui évaluent le risque marginal sont souvent appelées « études d'amélioration ». Les premières études ont indiqué, par exemple, que les chatbots de 2023 n'augmentaient pas significativement les risques de biosécurité par rapport aux technologies existantes : les participants ayant accès à Internet mais pas d'IA à usage général étaient en mesure d'obtenir des informations relatives aux armes biologiques à des taux similaires à ceux des participants ayant accès à l'IA (393) (voir [2.1.4. Attaques biologiques et chimiques](#) pour une discussion plus approfondie sur l'IA actuelle et le biorisque et [3.3. Identification et évaluation des risques](#) pour discussion sur les études de soulèvement et autres évaluations des risques. D'autre part, plusieurs études ont montré que la création de contenus NCII et CSAM a augmenté de manière significative en raison de la diffusion ouverte de modèles de génération d'images tels que Stable Diffusion (912*, 913) (voir [2.1.1. Dommages causés aux individus par des contenus falsifiés](#)). Il est important de tenir compte du risque marginal pour garantir que les interventions sont proportionnelles au risque posé (393, 911).

Toutefois, pour pouvoir procéder à une analyse du risque marginal, les entreprises ou les régulateurs doivent d'abord établir un seuil de risque tolérable stable (voir [3.1. Aperçu de la gestion des risques](#)) auquel le risque marginal peut être comparé afin d'éviter un scénario de « grenouille bouillante » (910). Même si une amélioration progressive de la capacité du modèle n'augmente que légèrement le risque marginal par rapport à la technologie préexistante, la superposition indéfinie de risques marginaux mineurs sur des risques marginaux mineurs pourrait entraîner une augmentation substantielle du risque au fil du temps et conduire par inadvertance à la mise sur le marché d'une technologie inacceptablement dangereuse. En revanche, l'amélioration de la résilience sociétale et le renforcement des capacités défensives pourraient contribuer à maintenir le risque marginal à un niveau bas même si les capacités et la « montée en puissance » du modèle progressent.

Français Une lacune clé en matière de données probantes concerne la question de savoir si les versions à pondération ouverte pour l'IA à usage général auront un impact positif ou négatif sur la concurrence et la concentration du marché. La publication publique des pondérations des modèles peut entraîner des impacts à la fois positifs et négatifs sur la concurrence, la concentration du marché et le contrôle (901, 910, 914, 915). À court terme, le partage de modèles à pondération ouverte protégé par une licence open source donne du pouvoir aux petits développeurs en aval en leur donnant accès à des technologies sophistiquées qu'ils ne pourraient pas se permettre autrement de créer, favorisant ainsi l'innovation et diversifiant le paysage applicatif. On estime qu'un investissement d'un milliard d'euros dans de nombreux types de logiciels open source (OSS) dans l'UE en 2018 a eu un impact économique de 65 à 95 milliards d'euros (916). On pourrait s'attendre à un impact similaire de l'IA à pondération ouverte publiée sous licence open source.

Cependant, cette démocratisation apparente de l'IA peut également jouer un rôle dans le renforcement de la domination et de la concentration du marché ([2.3.3. Concentration du marché et points de défaillance uniques](#)) parmi les principaux acteurs (914, 915). À plus long terme, les entreprises qui publient des modèles d'IA à usage général et ouverts voient souvent leurs cadres devenir des normes industrielles, façonnant l'orientation des développements futurs, comme cela devient rapidement le cas avec l'utilisation généralisée des modèles Llama dans les projets de développement ouverts et les applications industrielles. Ces entreprises peuvent ensuite facilement intégrer les avancées réalisées par la communauté (gratuitement) dans leurs propres offres, préservant ainsi leur avantage concurrentiel.

En outre, l'écosystème de développement open source plus large constitue un terrain de recrutement fertile.

Les entreprises peuvent ainsi identifier et attirer des professionnels qualifiés qui connaissent déjà leurs technologies (914). Il est probable que la diffusion de données ouvertes affectera différemment la concentration du marché à différents niveaux de l'écosystème de l'IA à usage général ; elle est plus susceptible d'accroître la concurrence et de réduire la concentration du marché dans le développement d'applications en aval, mais au niveau du développement de modèles en amont, la direction de l'effet est plus incertaine (750). Des recherches supplémentaires sont nécessaires pour clarifier la dynamique technique et économique en jeu.

Une fois que les poids des modèles sont disponibles pour téléchargement public, il n'existe aucun moyen de mettre en œuvre une restauration complète de toutes les copies existantes. Les plateformes d'hébergement Internet telles que GitHub et Hugging Face peuvent supprimer des modèles de leurs plateformes, ce qui rend difficile pour certains acteurs de trouver des copies téléchargeables et constitue une barrière suffisante pour de nombreux utilisateurs malveillants occasionnels à la recherche d'un moyen facile de causer du tort (917). Cependant, un acteur bien motivé serait toujours en mesure d'obtenir une copie de modèle ouvert malgré les inconvénients ; même les modèles de grande taille sont faciles à distribuer en ligne et hors ligne. Par exemple, les modèles de pointe tels que Llama-3.1-405B peuvent tenir sur une clé USB, ce qui souligne la difficulté de contrôler la distribution une fois que les modèles sont publiés ouvertement.

Les solutions techniques pour réduire les risques liés à la diffusion de poids ouverts sont encore émergentes et impliquent souvent des compromis importants par rapport aux avantages des modèles entièrement ouverts. Par exemple, les « modèles de récupération » permettent aux développeurs de partitionner les capacités « sûres » et « non sûres », ce qui permet potentiellement de diffuser ouvertement des parties non dangereuses d'un modèle tout en limitant les capacités dangereuses. Cependant, ces modèles sont confrontés à des défis tels que la rigidité contextuelle et nécessitent un accès aux données sources (918). D'autres techniques sont en cours de développement pour atténuer les abus en réduisant les performances du modèle lorsque les poids sont falsifiés (919, 920, 921, 922, 923, 924). Cependant, ces méthodes actuelles sont naissantes et souffrent de compromis majeurs en termes d'efficacité, de stabilité et de performances sur des tâches bénignes. L'établissement de repères et l'amélioration des techniques de « désapprentissage » inviolable restent un défi permanent, comme indiqué dans [3.4.1. Formation de modèles plus fiables.](#)

Il existe des approches d'atténuation des risques pour les modèles à pondération ouverte tout au long du cycle de vie de l'IA. Les stratégies d'atténuation des risques les plus robustes viseront à résoudre les problèmes potentiels à chaque étape (voir [3.1. Vue d'ensemble de la gestion des risques](#)), de la collecte de données et de la formation du modèle jusqu'aux réglages fins et aux mesures post-publication telles que la divulgation des vulnérabilités (informer les utilisateurs lorsqu'une faille du modèle a été trouvée) (910, 925). Même si les pondérations des modèles sont maintenues complètement fermées, ces approches d'atténuation permettent aux développeurs de planifier les fuites, car les atténuations des risques pour les modèles à pondération ouverte sont susceptibles d'être également utiles pour les modèles fermés divulgués. Par exemple, le modèle Llama 3.1 à 405 milliards de paramètres aurait été divulgué au public avant sa publication ouverte (926).

Pour les décideurs politiques qui travaillent à la réglementation de la diffusion des modèles, les principaux défis à venir sont les suivants :

- Rechercher des preuves de risque marginal dans les zones d'incertitude. Les décideurs politiques ont besoin de données solides
Analyses du risque marginal pour comprendre où l'ouverture introduit des risques significatifs et où elle n'en introduit pas. La plupart des recherches actuelles n'évaluent pas le risque marginal des modèles ouverts.
- Suivre et anticiper l'évolution des risques avec le développement technologique. En tant qu'IA
À mesure que les capacités progressent, les risques associés peuvent (parfois rapidement) augmenter (en raison de l'accès des adversaires à des modèles de capacités plus élevées) ou diminuer (en raison d'une IA plus fiable ou de meilleures défenses créées grâce à l'IA), nécessitant une évaluation et une adaptation continues des politiques (voir également [1.3. Capacités dans les années à venir](#)).
- Reconnaître que certaines interventions politiques ne peuvent pas être appliquées aux modèles ouverts en raison de limitations techniques. Par exemple, l'obligation d'apposer un filigrane sur les modèles linguistiques ne peut pas être appliquée aux modèles ouverts, car il existe des limitations techniques à la mise en œuvre de filigranes qui ne peuvent pas être supprimées.
- Connaître les interventions techniquement réalisables et la manière dont la diffusion ouverte affecte ces interventions. Par exemple, les interventions techniques ou les restrictions sur le réglage fin des modèles d'IA à usage général ne sont pas réalisables pour les modèles d'IA à pondération ouverte.
- Analyse de l'impact positif et négatif de la régulation de la diffusion des modèles. IA à pondération ouverte
les modèles présentent de solides avantages en termes de transparence, de concurrence et de concentration – du moins dans certaines parties de l'écosystème de l'IA.
- Équilibrer les risques marginaux avec les avantages marginaux. Il est important de développer des cadres pour prendre des décisions concernant la réglementation des modèles d'IA à pondération ouverte. Ces cadres sont susceptibles de dépendre du contexte et il n'existe pas de réponse unique et correcte : les différentes parties, institutions et gouvernements parviendront à des conclusions différentes en fonction de leurs priorités et des spécificités du modèle et du mécanisme de diffusion envisagés.

Pour les pratiques de gestion des risques liées aux modèles à pondération ouverte, voir :

- [3.1. Aperçu de la gestion des risques](#)
- [3.3. Identification et évaluation des risques](#)

3. Approches techniques de la gestion des risques

3.1. Aperçu de la gestion des risques

INFORMATIONS CLÉS

- La gestion des risques (identification et évaluation des risques, puis atténuation et suivi de ces risques) est un défi dans le contexte de l'IA à usage général. Bien que de nombreux cadres et pratiques soient en cours d'élaboration à l'échelle mondiale, des lacunes importantes subsistent en matière de validation, de normalisation et de mise en œuvre dans les différents secteurs et juridictions, en particulier pour identifier et atténuer les risques sans précédent.
- Le contexte de la gestion des risques liés à l'IA à usage général est particulièrement complexe en raison de l'évolution rapide de la technologie et de sa large applicabilité. Les pratiques traditionnelles de gestion des risques (telles que la sécurité dès la conception, les audits, la redondance et les dossiers de sécurité) fournissent une base, mais doivent être adaptées compte tenu de l'évolution rapide, de la large applicabilité et des effets d'interaction complexes de l'IA à usage général.
- Une approche de « sécurité du système » est utile pour gérer efficacement les risques généraux liés à l'IA. Cette approche applique à la fois des principes d'ingénierie et de gestion pour identifier et contrôler les dangers tout au long du cycle de vie d'un système. Pour l'IA à usage général, cela inclut la compréhension des interactions entre les composants matériels et logiciels, les structures organisationnelles et les facteurs humains.
- Une stratégie de « défense en profondeur » est apparue comme une approche technique de premier plan. La stratégie consistant à superposer plusieurs mesures de protection est courante dans des domaines tels que la sûreté nucléaire et le contrôle des maladies infectieuses. Elle est en cours d'adaptation aux systèmes d'IA à usage général tout au long de leur cycle de vie, avec des rôles différents pour les fournisseurs de données, les fournisseurs d'infrastructures, les développeurs et les utilisateurs.
- Les données actuelles mettent en évidence deux défis majeurs dans la gestion des risques liés à l'IA à usage général. Premièrement, il est difficile de hiérarchiser les risques en raison de l'incertitude quant à leur gravité et à leur probabilité d'occurrence. Deuxièmement, il peut être complexe de déterminer les rôles et responsabilités appropriés tout au long de la chaîne de valeur de l'IA et d'encourager une action efficace.

Définitions clés

- **Risque** : la combinaison de la probabilité et de la gravité d'un préjudice résultant du développement, du déploiement ou de l'utilisation de l'IA.
- **Danger** : tout événement ou activité susceptible de causer un préjudice, tel qu'une perte de vie, une blessure, une perturbation sociale ou des dommages environnementaux.
- **Gestion des risques** : processus systématique d'identification, d'évaluation, d'atténuation et surveillance des risques.
- **Défense en profondeur** : une stratégie qui comprend la superposition de plusieurs mesures d'atténuation des risques cas où aucune méthode unique existante ne peut assurer la sécurité.
- **Capacités** : l'éventail des tâches ou des fonctions qu'un système d'IA peut exécuter et comment il peut les exécuter avec compétence.

- Déploiement : processus de mise en œuvre de systèmes d'IA dans des applications, des produits ou des services du monde réel où ils peuvent répondre aux demandes et fonctionner dans un contexte plus large.
- Modalités : les types de données qu'un système d'IA peut recevoir avec compétence en entrée et produire en sortie, notamment du texte (langage ou code), des images, des vidéos et des actions robotiques.

Défis de la gestion des risques

Les premières étapes du processus de gestion des risques comprennent l'identification et l'évaluation des risques, qui sont difficiles et bénéficient d'une expertise diversifiée. Ces sujets sont abordés en détail dans [la section 3.3. Identification et évaluation des risques](#), mais il est essentiel de les garder à l'esprit pour la gestion globale des risques en raison de leurs défis uniques et de la manière dont ils influencent tous les éléments ultérieurs de la gestion des risques. Il est essentiel d'identifier et d'évaluer les risques de l'IA à usage général dès les premières étapes de conception et pas seulement après le développement d'un modèle. Cela peut être facilité par l'utilisation de taxonomies et de typologies de risques complètes, qui catégorisent et organisent un grand nombre de risques. Les étapes ultérieures du processus de gestion des risques, y compris la priorisation et l'atténuation, sont abordées tout au long de la [section 3. Approches techniques de la gestion des risques](#), ainsi que dans le [tableau des pratiques de gestion des risques ci-dessous](#).

L'identification et l'évaluation des risques demeurent un défi car l'IA à usage général peut être appliquée dans de nombreux domaines et contextes différents, et les capacités (et les risques associés) évoluent au fil du temps. L'IA peut présenter des risques très différents lorsqu'elle est appliquée, par exemple, dans le domaine des soins de santé (où la précision est essentielle) et dans l'écriture créative (où elle ne l'est pas). De plus, des études montrent que les performances des systèmes d'IA à usage général peuvent évoluer au fil du temps, car elles peuvent être considérablement améliorées par des mesures relativement simples sans recyclage coûteux. Pour y remédier, il peut être nécessaire de procéder à des évaluations des risques régulièrement mises à jour (77). Par exemple, le réglage fin des modèles (par exemple en leur fournissant de petites quantités de données de formation supplémentaires hautement organisées) peut améliorer considérablement leurs capacités dans des domaines spécifiques (927), avec des implications pour le risque abordées dans [la section 2.1.4. Attaques biologiques et chimiques](#). Certains risques peuvent ne pas être prévisibles et résulter d'interactions complexes entre les modèles, les personnes, les organisations et les systèmes sociaux et politiques (172).

Pour mieux éclairer les pratiques de gestion des risques, il est nécessaire de procéder à des évaluations qui se concentrent sur un ensemble plus large de risques liés à l'IA à usage général, et pas seulement sur les capacités, et d'améliorer les évaluations dans les langues, les cultures, les modalités et les cas d'utilisation. Comme indiqué dans [la section 3.3. Identification et évaluation des risques](#), des progrès récents ont été réalisés dans les méthodes d'évaluation, notamment le référentiel MLCommons AI Safety, qui mesure la sécurité des grands modèles linguistiques (LLM) en évaluant les réponses des modèles aux invites dans plusieurs catégories de dangers, notamment l'exploitation sexuelle des enfants, les armes indiscriminées, le suicide et l'automutilation (457). Le référentiel d'évaluation de la sécurité sociotechnique comprend de nombreux référentiels et méthodes d'évaluation supplémentaires qui peuvent aider les développeurs et les évaluateurs à évaluer les risques sociétaux liés aux LLM et à d'autres systèmes d'IA génératifs (928*). Cependant, l'espace manque d'une approche plus large de la science de l'évaluation. Les évaluations actuelles se concentrent en grande partie sur le modèle d'IA à usage général lui-même, passant sous silence les différentes conceptions de systèmes, les cas d'utilisation, les publics d'utilisateurs et d'autres facteurs contextuels qui influencent fortement la manière dont le risque peut se manifester. Beaucoup d'entre eux

Les modèles se concentrent sur les modalités textuelles et peuvent être moins pertinents pour d'autres modalités (comme les images et l'audio) ou pour les systèmes multimodaux (929*). Ils ont également du mal à évaluer avec précision les risques dans le monde entier car, par exemple, ils n'évaluent qu'en anglais en fonction d'un contexte culturel occidental, mais le modèle peut être conçu pour être un système multilingue (930*). L'amélioration des repères pour les modèles dans les langues à faibles ressources nécessite une collaboration entre les chercheurs, les locuteurs natifs et les partenaires communautaires tels que les militants linguistiques et les éducateurs (931).

Une large participation et un large engagement sont nécessaires pour évaluer et gérer les risques liés à l'IA à usage général ; cela ne peut pas être laissé entre les mains de la seule communauté scientifique. La gestion efficace des risques liés aux systèmes d'IA à usage général hautement performants nécessite l'implication de plusieurs groupes, notamment des experts de plusieurs domaines et des communautés concernées, pour identifier et évaluer les risques hautement prioritaires. Même les notions de « risque » et de « sécurité » sont controversées – par exemple, elles laissent ouverte la question de savoir qui est concerné par la sécurité – et leur évaluation nécessite l'implication de divers groupes d'experts et de populations concernées (537). Il est courant que les cadres de gestion des risques de l'IA recommandent des méthodes participatives, notamment l'engagement d'un large ensemble de groupes concernés tout au long du cycle de vie de l'IA ; les approches participatives peuvent être difficiles à mettre en œuvre face à diverses dynamiques de pouvoir (932).

Mécanismes et pratiques de gestion des risques

Il existe de nombreuses pratiques et mécanismes qui peuvent aider à gérer le large éventail de risques posés par l'IA à usage général. Certains d'entre eux sont référencés dans le tableau 3.1 ci-dessous ; ils sont abordés plus en détail dans la section 3. [Approches techniques de la gestion des risques](#).

Le tableau 3.1 ci-dessous inclut les pratiques de gestion des risques qui soutiennent cinq étapes (interconnectées) de la gestion des risques :

- Identification des risques : processus de recherche, de reconnaissance et de description des risques.
- Évaluation des risques : Le processus visant à comprendre la nature du risque et à déterminer le niveau de risque.
- Évaluation des risques : processus de comparaison des résultats de l'évaluation des risques avec les critères de risque pour déterminer si le risque et/ou son ampleur est/sont acceptable(s) ou tolérable(s). (Notez que le terme « évaluation » a plusieurs significations dans le contexte de l'IA et peut également faire référence aux modèles de test.)
- Atténuation des risques : hiérarchiser, évaluer et mettre en œuvre les contrôles/contre-mesures de réduction des risques appropriés recommandés par le processus de gestion des risques.
- Gouvernance des risques : processus par lequel l'évaluation, les décisions et les actions en matière de gestion des risques sont liées à la stratégie et aux objectifs de l'entreprise. La gouvernance des risques offre la transparence, la responsabilité et l'obligation de rendre des comptes qui permettent aux gestionnaires de gérer les risques de manière acceptable.

Il convient de noter que la terminologie exacte utilisée pour décrire les différentes étapes de la gestion des risques varie selon les principaux cadres et normes. Le tableau est destiné à être illustratif plutôt qu'exhaustif.

Risque Gestion Scène	Gestion des risques Pratique/Méthode	Explication	Domaines d'utilisation
Risque Identification	Taxonomie des risques	Un moyen de catégoriser et d'organiser les risques sur plusieurs dimensions	Il y a plusieurs taxonomies de risques bien connues pour l'IA (439, 933)
	Engagement avec Experts et communautés concernés	Les experts du domaine, les utilisateurs et les communautés concernées ont des informations uniques sur les risques probables	Des lignes directrices émergent pour une IA participative et inclusive (934)
	Méthode Delphi	Une technique de prise de décision de groupe qui utilise une série de questionnaires pour recueillir le consensus d'un panel d'experts	La méthode Delphi a été utilisée pour aider à identifier les principaux risques liés à l'IA (935)
	Modélisation des menaces	Un processus permettant d'identifier les menaces et les vulnérabilités d'un système	La modélisation des menaces est couramment utilisée pour soutenir la sécurité de l'IA dans la recherche et le développement de l'IA (936)
	Analyse de scénario	Élaborer des scénarios futurs plausibles et analyser la matérialisation des risques	L'analyse et la planification de scénarios sont largement utilisées dans tous les secteurs, y compris dans le secteur de l'énergie, et pour faire face aux incertitudes des systèmes électriques (937)
Risque Évaluation	Impact Évaluation	Un outil utilisé pour évaluer la impacts potentiels d'une technologie ou d'un projet	La loi européenne sur l'IA exige que les développeurs de systèmes d'IA à haut risque effectuent Impact sur les droits fondamentaux Évaluations (938)
	Audits	Un examen formel de la conformité d'une organisation aux normes, politiques et procédures, généralement effectué par une partie externe	L'audit de l'IA est un domaine en pleine croissance, mais s'appuie sur une longue histoire d'audit dans d'autres domaines, notamment la réglementation financière, environnementale et sanitaire (939)
	Équipe rouge	Un exercice dans lequel un groupe de personnes ou de systèmes automatisés se font passer pour un adversaire et attaquent les systèmes d'une organisation afin d'identifier les vulnérabilités	Le red-teaming est généralement utilisé dans le domaine de la cybersécurité, mais il est également devenu courant pour l'IA (940)

	Repères	Un test ou une mesure standardisée, souvent quantitative, utilisée pour évaluer et comparer les performances des systèmes d'IA sur un ensemble fixe de tâches conçues pour représenter une utilisation dans le monde réel	En 2023, l'IA avait atteint des performances de niveau humain sur de nombreux tests de référence importants en matière d'IA (731)
	Évaluation du modèle	Processus d'évaluation et mesurer les performances d'un système d'IA sur une tâche particulière	Il existe d'innombrables IA évaluations pour évaluer différentes capacités et risques, y compris pour la sécurité (941*)
	Analyse de sécurité	Permet de comprendre les dépendances entre les composants et le système dont ils font partie, afin d'anticiper comment les défaillances des composants pourraient entraîner des dangers au niveau du système	Cette approche est utilisée dans des domaines critiques pour la sécurité, par exemple pour anticiper et prévenir les accidents d'avion ou les accidents de réacteur nucléaire. effondrements
Risque Évaluation	Tolérance au risque	Le niveau de risque qu'une organisation est prête à prendre	En matière d'IA, les tolérances aux risques sont souvent laissées à la discrétion des entreprises d'IA, mais les régimes réglementaires peuvent aider à identifier les risques inacceptables qui sont légalement interdits (942)
	Seuils de risque	Limites quantitatives ou qualitatives qui distinguent les risques acceptables des risques inacceptables et déclenchent des actions de gestion des risques spécifiques lorsqu'elles sont dépassées	Seuils de risque pour L'IA à usage général est déterminée par une combinaison d'évaluations des capacités, de l'impact, du calcul, de la portée et d'autres facteurs (943, 944)
	Matrices de risque	Un outil visuel qui permet de hiérarchiser les risques en fonction de leur probabilité d'occurrence et de leur impact potentiel	Les matrices de risque sont utilisées dans de nombreuses industries et à de nombreuses fins, comme par exemple par les institutions financières pour évaluer le risque de crédit ou par les entreprises pour évaluer les éventuelles perturbations de leurs chaînes d'approvisionnement
	Méthode du nœud papillon	Une technique de visualisation quantitative et qualitative du risque, permettant une différenciation claire entre la gestion proactive et réactive du risque, destinée à aider à prévenir et à atténuer les risques d'accidents majeurs	Les compagnies pétrolières et les gouvernements nationaux utilisent la méthode du nœud papillon (945)

Risque Atténuation	La sécurité dès la conception	Une approche qui centre la sécurité des utilisateurs dans la conception et le développement des produits et services	Cette approche est courante dans les domaines critiques de l'ingénierie et de la sécurité, notamment l'aviation et l'énergie.
	« La sécurité de la Fonction prévue (SOTIF)	Une approche qui exige que les ingénieurs fournissent la preuve qu'un système est sûr lorsqu'il fonctionne comme prévu	Cette approche est utilisée dans de nombreux domaines de l'ingénierie, comme dans la construction. et essais de véhicules routiers (946)
	Défense en profondeur	L'idée selon laquelle plusieurs couches de défense indépendantes et superposées peuvent être mises en œuvre de telle sorte que si l'une d'elles échoue, les autres resteront efficaces	Un exemple vient du domaine des maladies infectieuses, où de multiples mesures préventives (par exemple, vaccins, masques, lavage des mains) peuvent être superposées pour réduire le risque global.
	Si-Alors Engagements	Un ensemble de protocoles et d'engagements techniques et organisationnels pour gérer les risques à différents niveaux à mesure que les modèles d'IA deviennent plus performants	Certaines entreprises qui développent une IA à usage général utilisent ces types d'engagements comme politiques de mise à l'échelle responsable ou cadres similaires (594*, 596*, 947*)
	Responsable Libérer et Déploiement Stratégies	Il existe un éventail de stratégies de publication et de déploiement pour l'IA, notamment des versions par étapes, un accès basé sur le cloud ou via API, des contrôles de sécurité de déploiement et des politiques d'utilisation acceptables	Il existe certaines pratiques émergentes de l'industrie qui se concentrent sur la publication et stratégies de déploiement pour l'IA à usage général (596*, 947*, 948)
	Dossiers de sécurité	Les dossiers de sécurité exigent que les développeurs démontrent la sécurité. Un dossier de sécurité est un argument structuré appuyé par des preuves qu'un système est acceptable pour fonctionner en toute sécurité dans un contexte particulier	Les cas de sécurité sont courants dans de nombreux secteurs, notamment la défense, l'aérospatiale et les chemins de fer (949)
Risque Gouvernance	Documentation	Il existe de nombreux bonnes pratiques de documentation, lignes directrices et exigences pour les systèmes d'IA afin de suivre, par exemple, les données de formation, la conception et la fonctionnalité du modèle, les cas d'utilisation prévus, les limitations et les risques	Les « cartes modèles » et les « cartes système » sont des exemples de normes de documentation d'IA importantes (34, 51*)
	Registre des risques	Un outil de gestion des risques qui sert de référentiel de tous les risques, de leur hiérarchisation, de leurs propriétaires et de leurs plans d'atténuation. Ils sont parfois utilisés pour se conformer à la réglementation	Les registres de risques sont un outil relativement standard utilisé dans de nombreux secteurs, notamment la cybersécurité (950) et récemment l'IA (933, 951*)

	<p>Dénonciateur Protection</p>	<p>Les lanceurs d'alerte peuvent jouer un rôle important en alertant les autorités des risques dangereux à Les entreprises d'IA en raison de la nature exclusive de nombreuses avancées de l'IA</p>	<p>Les incitations et les protections pour les lanceurs d'alerte sont devrait être un élément important de la gouvernance avancée des risques liés à l'IA (952)</p>
	<p>Rapport d'incident</p>	<p>Le processus de documentation et de partage systématique des cas dans lesquels le développement ou le déploiement L'IA a causé des dommages directs ou indirects nuit</p>	<p>Le signalement d'incidents est courant dans de nombreux domaines, des ressources humaines à cybersécurité. Il est également devenu plus courant pour l'IA (953)</p>
	<p>Gestion des risques Cadres</p>	<p>Des cadres organisationnels complets pour réduire les lacunes dans la couverture des risques et garantir que les diverses activités à risque (c'est-à-dire toutes celles mentionnées ci-dessus) sont structurées et alignées de manière cohérente, que les rôles et responsabilités en matière de risque sont clairement définis et que des freins et contrepoids sont en place pour éviter les silos et gérer les conflits d'intérêts.</p>	<p>Dans d'autres secteurs critiques pour la sécurité, les trois lignes de Cadre de défense – séparer la propriété des risques, la surveillance et l'audit – est largement utilisé et peut être utilement appliqué aux Entreprises d'IA (954, 955)</p>

Tableau 3.1 : Plusieurs pratiques et mécanismes, organisés en cinq étapes de gestion des risques, peuvent aider à gérer le large éventail de risques posés par l'IA à usage général.

Les mécanismes de documentation et de transparence institutionnelle, ainsi que les pratiques de partage d'informations, jouent un rôle important dans la gestion des risques liés à l'IA à usage général et facilitent le contrôle externe. Il est devenu courant de tester les modèles avant leur publication, notamment par le biais de red-teaming et d'analyses comparatives, et de publier les résultats dans une « fiche modèle » ou « fiche système » accompagnée de détails de base sur le modèle, notamment la manière dont il a été formé et ses limites (34, 51*). Une autre approche qui peut favoriser des niveaux plus élevés de transparence institutionnelle consiste à publier des rapports de transparence sur les modèles de base ou à rendre public un degré similaire de documentation (956).

D'autres éléments importants de la documentation et de la transparence comprennent la surveillance et le signalement des incidents (44*, 957*), par exemple via l'Initiative de partage des incidents liés à l'IA (953) ; et le partage d'informations, qui peut être facilité par des groupes industriels tels que le Frontier Model Forum, les gouvernements ou d'autres.

L'amélioration et la normalisation de la documentation favorisent un contrôle et une responsabilisation externes accrus (958).

La tolérance au risque et les seuils de risque sont des aspects particulièrement importants de la gestion des risques pour l'IA à usage général. Il n'est pas possible d'évaluer l'IA à usage général pour toutes les capacités possibles, de sorte que les organisations donnent la priorité à celles qui sont les plus susceptibles d'entraîner des résultats néfastes au-delà de leur tolérance au risque. La tolérance au risque est souvent laissée aux développeurs et aux déployeurs d'IA pour la déterminer eux-mêmes, mais les décideurs politiques peuvent aider à fournir des orientations et des restrictions concernant les risques inacceptables pour les individus et la société. Une pratique de plus en plus courante parmi les développeurs d'IA consiste à

Les seuils de capacités prédéfinis volontaires (594*, 947*) déterminent que lorsque les modèles présentent des capacités spécifiques (risquées), celles-ci doivent être satisfaites par des mesures d'atténuation spécifiques destinées à maintenir les risques à un niveau acceptable. Par exemple, une entreprise s'est engagée à mettre en œuvre une série de couches défensives (« défense en profondeur ») conçues pour empêcher toute utilisation abusive dès qu'un modèle est jugé « aider de manière significative des individus ou des groupes ayant une formation STEM de base à obtenir, produire ou déployer des armes CBRN [chimiques, biologiques, radiologiques et nucléaires] » (947*). Ces seuils de capacités peuvent avoir l'avantage d'être observables et mesurables dans une certaine mesure. Cependant, les capacités ne sont qu'un des multiples « indicateurs de risque clés » possibles et l'évaluation des capacités n'est pas une évaluation complète des risques. D'autres types de seuils pertinents pour l'IA à usage général incluent les seuils de risque, qui tentent d'estimer directement le niveau de risque (944) et les seuils de calcul, qui définissent des seuils en termes de ressources informatiques nécessaires pour entraîner un modèle (943). Cependant, des limites importantes demeurent. Les seuils de calcul en particulier ne sont pas un indicateur fiable du risque (170*), bien qu'ils présentent l'avantage d'être facilement mesurables, pertinents pour de nombreux risques différents et connus bien avant que les risques ne se matérialisent réellement. D'autres critères tels que le nombre d'utilisateurs privés ou professionnels, la gamme de modalités qu'une IA peut gérer, ainsi que la taille et la qualité des données d'entraînement, pourraient également jouer un rôle dans la définition des seuils de risque à l'avenir (959).

Français Les stratégies de diffusion et de déploiement de l'IA sont une pratique de gestion des risques supplémentaire qui peut être particulièrement utile pour l'IA à usage général. [2.4. Impact des modèles d'IA à usage général à pondération ouverte sur les risques liés à l'IA](#) examine comment la diffusion ouverte des pondérations des modèles affecte les risques. Certaines bonnes pratiques émergentes du secteur se concentrent sur les stratégies de diffusion et de déploiement de l'IA à usage général (948). Les stratégies de diffusion possibles incluent la diffusion du modèle par étapes pour tirer des enseignements des preuves du monde réel avant la diffusion complète, la fourniture d'un accès basé sur le cloud ou API (interface de programmation d'application) pour avoir une plus grande capacité à prévenir les abus, ou la mise en œuvre d'autres contrôles de sécurité du déploiement (44*). D'autres approches incluent l'utilisation de licences d'IA responsables et de politiques d'utilisation acceptables pour limiter potentiellement les abus (960).

Les pratiques de gestion des risques nécessitent un engagement de la part des dirigeants de l'organisation et des incitations organisationnelles alignées. La culture et la structure organisationnelles ont un impact sur l'efficacité des initiatives d'IA responsable et sur la gestion des risques liés à l'IA de nombreuses manières (961). Certains développeurs disposent de comités de décision internes qui délibèrent sur la manière de concevoir, de développer et d'examiner de nouveaux systèmes de manière sûre et responsable. Les comités de surveillance et de conseil, les fiduciaires ou les comités d'éthique de l'IA peuvent fournir des conseils utiles en matière de gestion des risques et de surveillance organisationnelle (962*, 963).

Leçons tirées d'autres domaines

Les stratégies de gestion des risques d'autres domaines peuvent être appliquées à l'IA à usage général. Les outils de gestion des risques courants dans d'autres secteurs critiques pour la sécurité, tels que la biosécurité et la sûreté nucléaire, comprennent des audits et des inspections planifiés, garantissant la traçabilité à l'aide d'une documentation standardisée, des mécanismes de défense redondants contre les risques et les défaillances critiques, des tampons de sécurité, des bandes de contrôle, des évaluations d'impact à long terme, ALARP (acronyme de « maintenir le risque aussi bas que raisonnablement possible »).

Les évaluations d'impact sur les droits de l'homme sont également utilisées dans de nombreux domaines pour évaluer les impacts sur les droits de l'homme de pratiques industrielles particulières (964), et sont très pertinentes pour les systèmes d'IA de toutes sortes (965). La prévision est une autre méthode de longue date, qui présente à la fois des avantages et des inconvénients (966), mais qui peut aider à éclairer les décisions à enjeux élevés concernant l'IA à usage général (928*, 967). Bien qu'il puisse être difficile de traduire les meilleures pratiques d'autres domaines en IA à usage général, il existe des conseils sur les moyens d'y parvenir (968, 969).

L'ingénierie de la sécurité et de la fiabilité est particulièrement pertinente. L'« ingénierie de la sécurité des systèmes » se concentre sur les interactions entre les différentes parties d'un système plus vaste (970) et souligne que les accidents peuvent survenir pour des raisons plus complexes que de simples défaillances de composants, des chaînes d'événements de défaillance ou des écarts par rapport aux attentes opérationnelles (971, 972). Dans le cas de l'IA, l'ingénierie de la sécurité des systèmes implique de prendre en compte tous les éléments constitutifs d'un système d'IA à usage général, ainsi que le contexte plus large dans lequel il fonctionne. La pratique de l'ingénierie de la sécurité a une longue histoire dans divers systèmes d'ingénierie critiques pour la sécurité, tels que le contrôle de vol des avions, les systèmes de contrôle des moteurs et le contrôle des réacteurs nucléaires. À un niveau élevé, l'ingénierie de la sécurité garantit qu'un système essentiel à la vie humaine fonctionne comme prévu et avec un minimum de dommages, même lorsque certains composants du système tombent en panne.

L'« ingénierie de fiabilité » a une portée plus large et aborde également les défaillances non critiques.

Ces approches offrent plusieurs techniques utiles pour l'évaluation des risques à usage général.

IA:

- La « sécurité dès la conception » (SbD) est une approche qui centre la sécurité des utilisateurs dans la conception et Développement de produits et de services. Pour les produits et services d'IA à usage général, cela peut prendre la forme d'une minimisation du contenu illégal, nuisible et dangereux dans les données de formation du modèle et d'une évaluation d'un large éventail de risques avant le déploiement.
- L'analyse de sécurité décrit les relations de cause à effet entre la fonctionnalité des composants individuels et le système global, de sorte que les défaillances des composants, qui peuvent entraîner des dangers au niveau du système (par exemple, des accidents d'avion ou des fusions du cœur d'un réacteur nucléaire), peuvent être anticipées et évitées dans la mesure du possible. Pour l'IA à usage général, cela pourrait signifier chercher à comprendre comment les pratiques de sécurité des données d'apprentissage d'un modèle peuvent influencer la sécurité du modèle global.
- Les approches de « sécurité de la fonction prévue » (SOTIF) exigent que les ingénieurs fournissent la preuve que le système est sûr lorsqu'il fonctionne comme prévu. SOTIF est particulièrement pertinent pour l'IA à usage général, car il prend en compte les scénarios dans lesquels un système peut fonctionner correctement mais présenter néanmoins un risque de sécurité en raison de circonstances imprévues.
- Certaines méthodes d'évaluation des risques, comme pour le secteur de l'énergie nucléaire, s'appuient sur modèles mathématiques conçus pour quantifier le risque en fonction de divers choix de conception et d'ingénierie, accompagnés de seuils de risque quantitatifs fixés par les régulateurs (973). Par exemple, certaines commissions de réglementation exigent que les exploitants de réacteurs nucléaires produisent des évaluations probabilistes des risques et veillent à ce que les risques estimés de certains événements soient conservés.

en dessous des seuils spécifiés. Bien que cela ne soit pas encore courant pour l'IA à usage général en raison des nombreux défis de quantification évoqués dans ce rapport, l'un des principaux avantages de cette approche est qu'elle permet à un organisme publiquement responsable de définir quel risque est considéré comme acceptable ou inacceptable, d'une manière accessible au public et aux experts externes.

Il est essentiel d'examiner attentivement les choix de conception effectués tout au long du cycle de vie de l'IA à usage général. L'approche « consciente du pipeline » pour atténuer les dommages causés par l'IA s'inspire de l'ingénierie de la sécurité et propose d'examiner attentivement de nombreux choix de conception effectués tout au long du cycle de vie de l'IA à usage général, de l'idéation et de la formulation des problèmes à la conception, au développement et au déploiement, à la fois en tant que composants individuels et en relation les uns avec les autres (974, 975). Des travaux supplémentaires sont nécessaires pour étendre ces idées de l'IA traditionnelle à l'IA à usage général. Par exemple, Assurance of Machine Learning for use in Autonomous Systems (AMLAS) fournit une méthodologie pour intégrer l'assurance de la sécurité dans le développement de composants d'apprentissage automatique et peut également être utile pour l'IA à usage général (976).

Les « dossiers de sécurité » pourraient fournir aux décideurs politiques un moyen utile d'explorer les dangers et les mesures d'atténuation des risques pour l'IA à usage général. Les développeurs de technologies critiques pour la sécurité telles que l'aviation, les dispositifs médicaux et les logiciels de défense sont tenus de réaliser des « dossiers de sécurité », qui imposent au développeur la charge de la preuve pour démontrer que son produit ne dépasse pas les seuils de risque maximum fixés par le régulateur (38, 949, 977, 978). Un dossier de sécurité est un argument structuré appuyé par des preuves, dans lequel le développeur identifie les dangers, modélise les scénarios de risque et évalue les mesures d'atténuation prises.

Par exemple, un dossier de sécurité pour une IA à usage général pourrait montrer qu'un système d'IA est incapable de provoquer des résultats inacceptables dans un contexte réaliste, par exemple même si le système est placé sur des serveurs non surveillés et dispose d'un accès à des ressources informatiques substantielles (978). Les dossiers de sécurité tirent parti de l'expertise technique du développeur de la technologie et peuvent être examinés par des tiers, mais nécessitent néanmoins que le régulateur (ou un tiers approprié) dispose de l'expertise technique et d'autres ressources pour les évaluer de manière appropriée. Une limitation possible est que les dossiers de sécurité peuvent ne traiter qu'un sous-ensemble de risques et de modèles de menaces, en laissant de côté les plus importants (979, 980). Une atténuation de cette limitation consiste à examiner les dossiers de sécurité parallèlement aux dossiers de risque produits par une équipe rouge d'experts tiers (978).

Le modèle de « défense en profondeur » est utile pour la gestion des risques généraux de l'IA. Il peut être judicieux de disposer de plusieurs niveaux de défense indépendants et superposés contre les risques, de sorte que si l'un d'eux échoue, les autres restent efficaces. C'est ce que l'on appelle parfois le « modèle de défense en fromage suisse ».

(981). Un exemple de l'efficacité du modèle de défense en profondeur est la gamme de mesures préventives déployées pour prévenir les maladies infectieuses : les vaccins, les masques et le lavage des mains, entre autres mesures, peuvent réduire considérablement le risque d'infection en combinaison, même si aucune de ces méthodes n'est efficace à 100 % à elle seule (981). Pour l'IA à usage général, la défense en profondeur comprendra des contrôles qui ne concernent pas le modèle d'IA lui-même, mais l'écosystème plus large, tels que les contrôles des données d'entraînement (par exemple certaines séquences d'ADN) et les contrôles des matériaux nécessaires à l'exécution d'une attaque (par exemple, l'équipement et les réactifs). Il est également important de se rappeler que des méthodes comme la défense en profondeur ne sont probablement pas suffisantes à elles seules, car elles se concentrent sur la prévention des accidents, des risques de dysfonctionnement (voir [2.2. Risques de dysfonctionnement](#)) et

Les risques d'utilisation malveillante sont réduits (voir [2.1. Risques liés à une utilisation malveillante](#)), mais ne sont généralement pas suffisants pour gérer des risques systémiques plus complexes (voir [2.3. Risques systémiques](#)).

Lacunes et opportunités

Les principales lacunes en matière de données probantes concernant la gestion des risques liés à l'IA à usage général concernent l'ampleur des risques et la mesure dans laquelle les différents mécanismes peuvent réellement limiter et atténuer les risques dans des contextes réels. Il n'existe pas toujours de consensus scientifique sur la probabilité ou la gravité des risques liés aux systèmes d'IA à usage général, ce qui rend difficile pour les décideurs politiques de savoir s'ils doivent être prioritaires et comment. Par exemple, la manière de gérer le risque d'utilisation abusive dépendra du niveau de compétence des acteurs de la menace dans des contextes réels préoccupants. En outre, la plupart des efforts de gestion des risques décrits ci-dessus ne sont pas encore validés, normalisés ou largement utilisés. Les efforts de gestion des risques varient selon les principales entreprises d'IA et les incitations peuvent ne pas être bien alignées pour encourager une évaluation et une gestion approfondies (982). Bien qu'il existe quelques mesures d'atténuation des risques qui sont perçues comme les plus efficaces par les experts pour réduire les risques systémiques liés à l'IA à usage général (983), l'efficacité des mécanismes de gestion des risques de l'IA à usage général est toujours en cours d'évaluation et les décideurs politiques devraient rechercher davantage de preuves issues d'applications concrètes.

Pour les décideurs politiques travaillant sur la gestion des risques liés à l'IA à usage général, les principaux défis consistent notamment à savoir comment hiérarchiser les nombreux risques posés par l'IA à usage général et à savoir qui est le mieux placé pour les atténuer. Les orientations en matière de gestion des risques recommandent souvent de donner la priorité aux préoccupations à forte probabilité ou à fort impact, notamment les cas où des impacts négatifs importants sont imminents ou déjà en cours, ou lorsque des risques catastrophiques pourraient être présents (887). Cependant, il n'est pas toujours évident de savoir quels sont les risques les plus probables ou les plus impactants. De plus, la gestion des risques implique nécessairement différents acteurs à différentes étapes de la chaîne de valeur de l'IA, notamment les fournisseurs de données et de cloud, les développeurs de modèles et les plateformes d'hébergement de modèles, chacun d'entre eux ayant des opportunités et des responsabilités uniques pour évaluer et gérer les risques. Les décideurs politiques doivent mieux comprendre en quoi les responsabilités des différents acteurs diffèrent et comment les incitations politiques peuvent soutenir diverses activités de gestion des risques (925).

3.2. Défis généraux en matière de gestion des risques et d'élaboration des politiques

3.2.1. Défis techniques pour la gestion des risques et l'élaboration des politiques

INFORMATIONS CLÉS

Plusieurs propriétés techniques de l'IA à usage général rendent difficile l'atténuation des risques associés à de nombreux risques liés à l'IA à usage général :

- A. Les agents d'IA polyvalents autonomes peuvent accroître les risques : les développeurs d'IA déploient de gros efforts pour créer et déployer des systèmes d'IA polyvalents capables d'agir et de planifier plus efficacement pour atteindre leurs objectifs. Ces agents ne sont pas bien compris, mais nécessitent une attention particulière de la part des décideurs politiques. Ils pourraient permettre des utilisations malveillantes et des risques de dysfonctionnements, tels que le manque de fiabilité et la perte de contrôle humain, en permettant des applications plus répandues avec moins de surveillance humaine.
- B. L'étendue des cas d'utilisation complique l'assurance de la sécurité : les systèmes d'IA à usage général sont étant utilisés pour de nombreuses tâches (souvent imprévues) dans de nombreux contextes, ce qui rend difficile de garantir leur sécurité dans tous les cas d'utilisation pertinents et permet potentiellement aux entreprises d'adapter leurs systèmes pour contourner les réglementations.
- C. Les développeurs d'IA à usage général comprennent peu le fonctionnement interne de leurs modèles : Malgré les progrès récents, les développeurs et les scientifiques ne peuvent pas encore expliquer pourquoi ces modèles génèrent un résultat donné, ni quelle fonction remplissent la plupart de leurs composants internes. Cela complique la garantie de sécurité et il n'est pas encore possible de fournir des garanties de sécurité même approximatives.
- D. Les comportements nuisibles, y compris les comportements non intentionnels orientés vers un objectif, persistent : Malgré les progrès progressifs réalisés dans l'identification et la suppression des comportements et des capacités nuisibles des systèmes d'IA à usage général, les développeurs peinent à les empêcher de présenter des comportements manifestement nuisibles dans des circonstances prévisibles, même ceux qui sont bien connus, comme la fourniture d'instructions pour des activités criminelles. En outre, les systèmes d'IA à usage général peuvent agir conformément à des objectifs imprévus qui peuvent être difficiles à prévoir et à atténuer.
- E. Un « déficit d'évaluation » en matière de sécurité persiste : malgré les progrès en cours, les méthodes actuelles d'évaluation et d'analyse des risques pour les systèmes d'IA à usage général sont immatures. Même si un modèle passe avec succès les évaluations de risques actuelles, il peut être dangereux. Pour développer les évaluations nécessaires à temps pour répondre aux engagements de gouvernance existants, des efforts, du temps, des ressources et un accès importants sont nécessaires.
- F. Les failles du système peuvent avoir un impact mondial rapide : lorsqu'un seul système d'IA polyvalent est largement utilisé dans plusieurs secteurs, des problèmes ou des comportements nuisibles peuvent affecter de nombreux utilisateurs simultanément. Ces impacts peuvent se manifester soudainement, par exemple lors de mises à jour de modèles ou d'une première version, et peuvent être pratiquement irréversibles.

Définitions clés

- **Agent IA** : une IA à usage général qui peut élaborer des plans pour atteindre des objectifs, effectuer de manière adaptative des tâches impliquant plusieurs étapes et des résultats incertains en cours de route, et interagir avec son environnement (par exemple en créant des fichiers, en effectuant des actions sur le Web ou en déléguant des tâches à d'autres agents) avec peu ou pas de surveillance humaine.
- **Déploiement** : processus de mise en œuvre de systèmes d'IA dans des applications, des produits ou des services du monde réel où ils peuvent répondre aux demandes et fonctionner dans un contexte plus large.
- **Évaluations** : Évaluations systématiques des performances, des capacités et des performances d'un système d'IA. vulnérabilités ou impacts potentiels. Les évaluations peuvent inclure des analyses comparatives, des équipes rouges et des audits et peuvent être menées avant et après le déploiement du modèle.
- **Réglage fin** : processus d'adaptation d'un modèle d'IA pré-entraîné à une tâche spécifique ou de le rendre plus utile en général en l'entraînant sur des données supplémentaires.
- **Généralisation erronée des objectifs** : situation dans laquelle un système d'IA suit correctement un objectif dans son environnement de formation, mais l'applique de manière inattendue lorsqu'il fonctionne dans un environnement différent.
- **Recherche d'interprétabilité** : l'étude du fonctionnement interne des modèles d'IA à usage général et le développement de méthodes pour rendre cela compréhensible pour les humains.
- **Jailbreaking** : Génération et envoi d'invites conçues pour contourner les garde-fous et amener un système d'IA à produire du contenu nuisible, tel que des instructions pour la construction d'armes.
- **Domaines ouverts** : environnements dans lesquels les systèmes d'IA peuvent être déployés présentent un très grand nombre de scénarios possibles. Dans les domaines ouverts, les développeurs ne peuvent généralement pas anticiper et tester toutes les manières possibles dont un système d'IA pourrait être utilisé.
- **Modèle à pondération ouverte** : un modèle d'IA dont les pondérations sont disponibles publiquement en téléchargement, comme Llama ou Stable Diffusion. Les modèles à pondération ouverte peuvent être, mais ne sont pas nécessairement, open source.
- **Pondérations** : paramètres du modèle qui représentent la force de la connexion entre les nœuds d'un réseau neuronal. Les pondérations jouent un rôle important dans la détermination de la sortie d'un modèle en réponse à une entrée donnée et sont mises à jour de manière itérative pendant l'entraînement du modèle pour améliorer son performance.

Cette section couvre six défis techniques généraux qui peuvent rendre la gestion des risques et l'élaboration des politiques plus difficiles pour un large éventail de risques liés à l'IA à usage général (voir la figure 3.1).

A. Les agents d'IA autonomes à usage général peuvent augmenter les risques : agents d'IA à usage général –

Les systèmes capables de planifier et d'agir dans le monde avec peu ou pas d'intervention humaine augmentent les risques de dysfonctionnements et d'utilisation malveillante. Aujourd'hui, les systèmes d'IA à usage général sont principalement utilisés comme outils par les humains. Par exemple, un chatbot peut écrire du code informatique, mais un humain exécute, débogue et intègre le code dans un projet logiciel plus vaste. Cependant, les chercheurs et les développeurs déploient de gros efforts pour concevoir des agents d'IA à usage général – des systèmes capables d'agir et de planifier de manière autonome en contrôlant des ordinateurs, des interfaces de programmation, des outils robotiques et en déléguant à d'autres systèmes d'IA (18, 55, 316*, 984, 985, 986*, 987, 988, 989, 990, 991*, 992). Ces systèmes sont également parfois appelés « agents autonomes » ou « IA autonome ». Les chercheurs et les développeurs construisent

agents pour une variété de domaines, y compris la navigation sur le Web (85*), la recherche en chimie et en IA (22*, 121*, 402), l'ingénierie logicielle (122, 259), la cybercriminalité (127), l'utilisation générale de l'ordinateur (993, 994*, 995) et le contrôle des robots (19*).

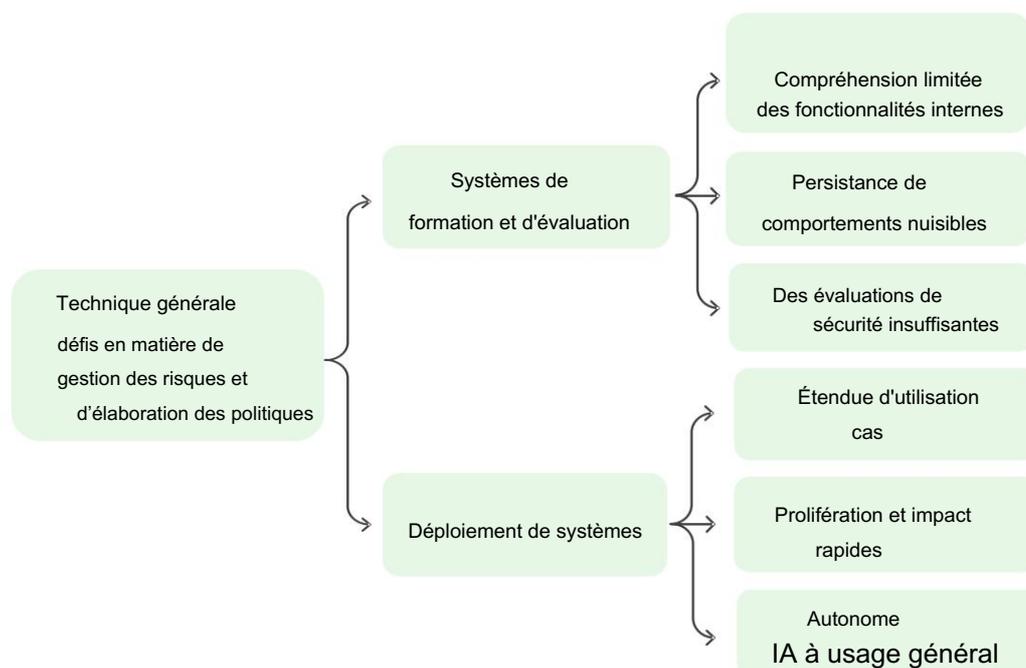


Figure 3.1 : Les défis techniques liés à la gestion des risques généraux liés à l'IA peuvent être divisés en deux types : les défis liés à la formation et à l'évaluation des systèmes, et les défis liés à leur déploiement. Cette section aborde six grands défis qui s'appliquent à de nombreux risques. Source : International AI Safety Report.

Les systèmes d'IA à usage général augmentent les risques en réduisant l'implication et la surveillance humaines.

L'objectif principal des agents d'IA à usage général est de réduire le besoin d'intervention et de surveillance humaines, ce qui permet des applications beaucoup plus rapides et moins chères. Cela est économiquement intéressant, et des produits d'IA de plus en plus agentsiques sont rapidement développés et déployés. Cependant, une délégation accrue aux agents d'IA réduit la surveillance humaine et peut augmenter le risque d'accidents (996) (voir 2.2.1. [Problèmes de fiabilité](#)). Dans le même temps, les agents peuvent être particulièrement vulnérables aux attaques d'acteurs malveillants (997), par exemple en « détournant » un agent en plaçant des instructions à des endroits où l'agent les rencontrera (998). Les agents d'IA peuvent également automatiser certains flux de travail pour des utilisations malveillantes telles que les escroqueries, le piratage et le développement d'armes (127, 358, 999, 1000, 1001*) (voir 2.1. [Risques liés à une utilisation malveillante](#) pour plus d'exemples). Les agents d'IA pourraient également contribuer de manière unique aux risques de perte de contrôle humain si leurs capacités progressent de manière significative (voir 2.2.3. [Perte de contrôle](#)) (316*, 1002).

En outre, les chercheurs ont fait valoir qu'il serait difficile, voire impossible, de garantir la sécurité des agents avancés en s'appuyant sur des tests, si ces agents peuvent établir des plans à long terme et distinguer les conditions de test des conditions réelles (1003).

Les agents d'IA polyvalents peuvent effectuer des tâches utiles de manière autonome, mais présentent actuellement une fiabilité limitée, en particulier pour les tâches complexes. Les systèmes d'IA polyvalents actuels sont capables d'exécuter de manière autonome de nombreuses tâches simples (par exemple, écrire de courts extraits de code), mais ils ont du mal à effectuer des tâches plus complexes (par exemple, écrire des bibliothèques de codes entières) (122, 593, 600, 1004).

Ils sont particulièrement peu fiables pour effectuer des tâches impliquant de nombreuses étapes (1005). Parallèlement, les agents d'IA polyvalents déployés pour accomplir des tâches à long terme peuvent être particulièrement vulnérables à la manipulation par des acteurs malveillants (997). Les capacités des agents actuels et futurs sont décrites plus en détail dans [1.2. Capacités actuelles](#) et [1.3. Capacités dans les années à venir](#).

Les capacités des agents d'IA polyvalents progressent rapidement, et la compréhension de leurs capacités futures constitue une lacune essentielle en matière de données probantes. Les agents d'IA polyvalents deviennent rapidement plus performants. Par exemple, « SWE-Bench » est une référence (métrique) populaire utilisée pour évaluer les capacités des systèmes d'IA agentiques pour les tâches d'ingénierie logicielle telles que la recherche et la correction de bugs (122). Depuis le rapport intermédiaire (mai 2024), les performances des meilleurs modèles sur SWE-Bench sont passées de 26 % à 42 % (122), les 19 meilleures soumissions ayant toutes eu lieu après mai 2024. Cela représente un progrès spectaculaire par rapport à octobre 2023, lorsque le meilleur modèle n'atteignait que 2 %. Parallèlement, l'introduction récente de o1 (2*) marque un bond en avant dans les capacités de raisonnement et de résolution de problèmes des systèmes d'IA à usage général. Ces améliorations de performances sont dues à une combinaison d'avancées. Tout d'abord, à mesure que les modèles d'IA à usage général sous-jacents à ces agents deviennent plus performants, les capacités cognitives des agents s'améliorent. Deuxièmement, ces agents sont développés avec des méthodes de formation et de planification de plus en plus avancées. Par exemple, AlphaProof, un système d'IA « neurosymbolique » à usage général qui combine des réseaux neuronaux avec des techniques de planification avancées, a obtenu des résultats dignes d'une médaille d'argent aux questions de l'Olympiade internationale de mathématiques 2024 (187*). Cependant, en raison du rythme rapide des progrès dans ce domaine et du fait que de nombreux agents sont propriétaires, la compréhension publique des méthodes de pointe actuelles est limitée. Au cours des mois et des années à venir, le développement d'agents plus avancés exige une attention particulière de la part des décideurs politiques.

B. L'étendue des cas d'utilisation complique l'assurance de la sécurité : les systèmes d'IA à usage général peuvent être appliqués dans de nombreux contextes imprévus, ce qui rend difficile de tester et de garantir leur fiabilité dans tous les cas d'utilisation réalistes. Les entrées et sorties des systèmes d'IA à usage général sont souvent ouvertes, comme la génération de texte ou d'images à forme libre où les utilisateurs peuvent saisir n'importe quelle invite. Il n'est pas possible d'étudier les impacts diffus et en aval d'un système dans un laboratoire de pré-déploiement. Il est donc difficile de fournir de solides garanties de sécurité, car il est impossible de tester de manière exhaustive un système dans tous les contextes d'utilisation pertinents. Par exemple, il existe des milliers de langues parlées par les humains, ce qui rend très difficile de garantir de manière exhaustive la sécurité des modèles linguistiques dans toutes les langues. Depuis la publication du rapport intermédiaire (mai 2024), les systèmes d'IA à usage général capables de traiter plusieurs types de données (par exemple, du texte, des images et de l'audio) sont devenus de plus en plus courants (1006). Cela élargit considérablement l'ensemble des contextes qui pourraient amener le système à se comporter de manière nuisible (1007). Les entreprises d'IA peuvent facilement réorienter les capacités de leurs systèmes entre différentes applications et solutions de contournement juridiques, ce qui pose des défis pour les approches d'intervention ciblées comme on l'a vu historiquement sur les marchés financiers (1008).

C. Les développeurs d'IA à usage général comprennent peu le fonctionnement interne de leurs modèles.

Une caractéristique clé des modèles d'IA à usage général est que leurs capacités sont principalement obtenues par l'apprentissage plutôt que par une conception descendante : un algorithme automatique ajuste des milliards de nombres

('paramètres') des millions de fois jusqu'à ce que la sortie du modèle corresponde aux données d'entraînement. En conséquence, la compréhension actuelle des modèles d'IA à usage général est plus analogue à celle des cerveaux en croissance ou des cellules biologiques qu'à celle des avions ou des centrales électriques. Les scientifiques et les développeurs d'IA n'ont qu'une capacité minimale à expliquer pourquoi ces modèles ont pris une décision donnée plutôt qu'une autre, et comment leurs capacités découlent de leurs composants mathématiques internes connus. Cela contraste, par exemple, avec les systèmes logiciels complexes tels que les moteurs de recherche Web, où les développeurs peuvent expliquer la fonction de composants individuels (tels que des lignes et des fichiers de code) et peuvent également rechercher pourquoi le système a trouvé un résultat particulier. Les techniques actuelles d'« interprétabilité » pour expliquer les structures internes des modèles d'IA à usage général ne sont pas fiables et nécessitent des hypothèses simplificatrices majeures (1009, 1010*, 1011*, 1012, 1013*). En pratique, les techniques d'interprétation du fonctionnement interne des réseaux neuronaux peuvent être trompeuses (466, 1014, 1015*, 1016, 1017, 1018, 1019) et peuvent échouer aux contrôles de cohérence ou s'avérer inutiles dans les utilisations en aval (1020, 1021, 1022, 1023, 1024, 1025*). Par exemple, l'un des objectifs de la recherche sur l'interprétabilité est d'aider les chercheurs à comprendre suffisamment bien les modèles pour modifier leurs comportements en modifiant leurs pondérations. Cependant, les outils d'interprétabilité de pointe ne se sont pas encore révélés utiles et fiables pour cela (1026*). Comme indiqué dans [3.4.1. Entraîner des modèles plus fiables](#),

Ces méthodes de recherche sont en cours d'amélioration et de nouveaux développements pourraient apporter de nouvelles perspectives. Cependant, en raison de la manière dont les modèles d'apprentissage profond représentent les informations à travers les neurones de manière hautement distribuée (1027, 1028), il n'est pas certain que l'interprétation des structures internes des modèles d'IA à usage général puisse offrir des garanties de sécurité. En d'autres termes, les systèmes d'IA à usage général modernes peuvent être trop complexes pour pouvoir garantir leurs performances. À l'heure actuelle, les informaticiens ne sont pas en mesure de donner des garanties du type « le système X ne fera pas Y » (41).

Néanmoins, une compréhension plus approfondie du fonctionnement interne des modèles pourrait être utile à de nombreux égards ([voir 3.4.2. Suivi et intervention](#) et [3.4.1. Formation de modèles plus fiables](#)).

D. Les comportements nuisibles, y compris les comportements non intentionnels orientés vers un objectif, persistent : veiller à ce que Il est difficile de démontrer que les systèmes d'IA à usage général agissent conformément aux objectifs, aux comportements et aux capacités prévus par leurs développeurs et utilisateurs. Bien que les systèmes d'IA à usage général puissent exceller dans l'apprentissage de ce qu'on leur « dit » de faire, leur comportement peut ne pas nécessairement correspondre à ce que leurs concepteurs avaient prévu (607, 1029, 1030, 1031). Même des différences subtiles entre les objectifs d'un concepteur et les objectifs assignés à un système peuvent conduire à des échecs inattendus. Par exemple, les chatbots d'IA à usage général sont souvent formés pour produire du texte qui sera évalué positivement par les évaluateurs, mais l'approbation de l'utilisateur est un indicateur imparfait des avantages pour l'utilisateur. En conséquence, plusieurs chatbots largement utilisés ont affiché un comportement « flagorneur » ou activement trompeur, faisant des déclarations que les utilisateurs approuvent, qu'elles soient vraies ou non (98, 317, 522, 608). Par exemple, les modèles de langage d'IA à usage général sont connus pour avoir une forte tendance à être d'accord avec les opinions qu'un utilisateur exprime dans les chats (98). Même lorsqu'un système d'IA à usage général reçoit un retour d'information correct pendant la formation, il peut toujours développer une solution qui ne se généralise pas bien lorsqu'elle est appliquée à de nouvelles situations une fois déployée (« généralisation erronée des objectifs ») (616, 1032, 1033). Par exemple, certains chercheurs ont découvert que la formation à la sécurité des modèles linguistiques peut être inefficace si le modèle est sollicité dans une langue qui était sous-représentée dans ses données de formation (1034). Depuis la publication du rapport intermédiaire (mai 2024), les chercheurs ont démontré des exemples de comportements indésirables axés sur les objectifs de la part de systèmes d'IA à usage général. Il s'agit notamment de tentatives de réécriture de leurs propres objectifs (599*).

Malgré les efforts déployés pour diagnostiquer et déboguer les problèmes, les développeurs n'ont pas toujours été en mesure d'empêcher même les comportements bien connus et manifestement nuisibles des systèmes d'IA à usage général dans des circonstances prévisibles. Empiriquement, les systèmes d'IA à usage général de pointe ont présenté une variété de comportements nuisibles et souvent inattendus après leur déploiement (41, 1035, 1036). Ces dangers incluent les systèmes d'IA à usage général qui aident les utilisateurs malveillants dans des tâches manifestement nuisibles (127, 319, 1037, 1038, 1039, 1040, 1041) ; la fuite d'informations privées ou protégées par des droits d'auteur (1042, 1043, 1044*, 1045, 1046, 1047) ; la génération de contenu haineux (1048, 1049) ; présentant des préjugés sociaux et politiques (183, 438, 491, 511, 560, 561, 562, 563, 564, 565) ; favorisant les préjugés des utilisateurs (98) ; et hallucinant un contenu inexact (101, 102*, 104, 461, 1050, 1051*). Entre-temps, les utilisateurs ont toujours pu contourner les mesures de protection des modèles d'IA à usage général de pointe avec une relative facilité grâce à des invites (« jailbreaks ») (39, 155, 460, 904*, 1052, 1053, 1054, 1055, 1056*, 1057, 1058, 1059, 1060, 1061, 1062, 1063*) ou à de simples modifications de modèle (906, 1064, 1065, 1066, 1067, 1068, 1069, 1070, 1071, 1072, 1073, 1074, 1075, 1076, 1077, 1078, 1079, 1080). Depuis la publication du rapport intermédiaire (mai 2024), certains chercheurs ont également constaté que même lorsque les systèmes de chat refusent en toute sécurité les demandes nuisibles, ils peuvent toujours se comporter de manière nuisible lorsqu'ils sont utilisés pour fonctionner comme des agents (1000, 1001*).

Les chercheurs développent continuellement de nouvelles techniques de défense contre ces attaques, mais ils développent également des attaques plus puissantes qui surmontent généralement les défenses existantes (voir [3.4.1 Former des modèles plus fiables](#)).

Français Les systèmes d'IA à usage général acquièrent et conservent parfois des capacités nuisibles même lorsqu'ils sont explicitement réglés pour ne pas le faire (41, 1069). Si les techniques actuelles sont efficaces pour supprimer les comportements nuisibles des systèmes d'IA à usage général, ces capacités nuisibles peuvent réapparaître et réapparaissent effectivement à partir d'anomalies, d'entrées d'utilisateurs malveillants et de modifications des modèles. Par exemple, le réglage fin de GPT-3.5 sur seulement dix exemples de texte nuisible peut annuler ses mesures de protection et permettre de susciter un comportement nuisible (1064). La difficulté de rendre les systèmes d'IA à usage général totalement résistants aux modes de défaillance manifestes a conduit certains chercheurs à se demander s'il est possible de rendre les approches de développement actuelles robustes à de tels modes de défaillance (1081, 1082). Voir [2.1. Risques d'utilisation malveillante](#) pour une discussion plus approfondie des capacités nuisibles dans les modèles d'IA, [2.4. Impact des modèles d'IA à usage général à pondération ouverte sur les risques de l'IA](#) pour une discussion des avantages et des risques de la publication de modèles avec des capacités à la fois nuisibles et bénéfiques pour le téléchargement public, et [3.4.1 Former des modèles plus fiables](#) pour une discussion sur les méthodes permettant de désapprendre les capacités nuisibles.

E. Un « déficit d'évaluation » en matière de sécurité persiste : les évaluations de sécurité actuelles ne sont pas suffisamment approfondies pour répondre aux cadres de gouvernance et aux engagements existants des entreprises. Les développeurs et les régulateurs proposent de plus en plus de cadres de gestion des risques qui s'appuient sur des évaluations de haute qualité des systèmes d'IA à usage général. L'objectif des évaluations est d'identifier les risques afin qu'ils puissent être traités ou surveillés. Cependant, la science de l'évaluation des systèmes d'IA à usage général et de la prévision de leurs impacts en aval est immature. Même lorsque les systèmes d'IA à usage général sont évalués avant le déploiement, de nouveaux modes de défaillance sont souvent rapidement découverts après le déploiement (1055). Par exemple, les utilisateurs ont trouvé des méthodes pour subvertir le réglage fin de la sécurité d'o1 dans les jours suivant sa sortie, et certains chercheurs ont rendu public des travaux sur une méthode permettant de jailbreaker de manière fiable le modèle seulement trois semaines après la sortie du modèle (1083).

Les risques en aval sont un domaine en pleine expansion. Cependant, la vaste portée des risques potentiels (933), les limites des techniques d'analyse comparative (178, 1084, 1085), le manque d'accès complet aux systèmes (1086) et la difficulté d'évaluer les impacts sociétaux en aval (928*, 930*, 933) rendent difficile la réalisation d'évaluations de haute qualité. [3.3. L'identification et l'évaluation des risques](#) approfondiront les méthodes d'évaluation des risques et les approches plus larges d'évaluation des risques.

F. Les failles du système peuvent avoir un impact mondial rapide : étant donné que les systèmes d'IA à usage général peuvent être partagés rapidement et déployés dans de nombreux secteurs (comme d'autres logiciels), un système nuisible peut rapidement avoir un impact mondial et parfois irréversible. Un petit nombre de modèles d'IA à usage général, propriétaires ou disponibles gratuitement, à pondération ouverte, touchent actuellement plusieurs millions d'utilisateurs (voir [2.3.3. Risques de concentration du marché et points de défaillance uniques](#)). Les modèles propriétaires et à pondération ouverte peuvent donc avoir des impacts rapides et mondiaux, bien que de différentes manières (911). Un facteur de risque pour les modèles à pondération ouverte est qu'il n'existe aucun moyen pratique de revenir en arrière si l'on découvre ultérieurement qu'un modèle présente des défauts ou des capacités qui permettent une utilisation malveillante (902) (voir [2.4. Impact des modèles d'IA à usage général à pondération ouverte sur les risques liés à l'IA, 2.1. Risques liés à une utilisation malveillante](#)). Cependant, l'un des avantages de la publication ouverte des pondérations des modèles et d'autres composants du modèle tels que le code et les données d'entraînement est qu'elle permet également à un nombre beaucoup plus grand et plus diversifié de praticiens de découvrir des failles, ce qui peut améliorer la compréhension des risques et des atténuations possibles (911). Les développeurs ou d'autres personnes peuvent alors réparer les défauts et proposer des versions nouvelles et améliorées du système. Cela ne peut pas empêcher une utilisation malveillante délibérée (902, 1075), ce qui pourrait être préoccupant si un système présente un risque supplémentaire (« risque marginal ») par rapport à l'utilisation d'alternatives (comme la recherche sur Internet). Tous ces facteurs sont pertinents pour la possibilité spécifique d'impacts rapides, généralisés et irréversibles des modèles d'IA à usage général. Cependant, même lorsque les composants du modèle ne sont pas rendus accessibles au public, les fonctionnalités du modèle atteignent toujours une large base d'utilisateurs dans de nombreux secteurs. Par exemple, dans les deux mois suivant son lancement, le système entièrement fermé ChatGPT comptait plus de 100 millions d'utilisateurs (1087).

3.2.2. Défis sociétaux en matière de gestion des risques et d'élaboration des politiques

INFORMATIONS CLÉS

Plusieurs facteurs économiques, politiques et autres facteurs contextuels rendent difficile l'atténuation des risques liés à de nombreux risques associés à l'IA à usage général :

- A. À mesure que l'IA à usage général progresse rapidement, l'évaluation des risques, l'atténuation des risques, la gouvernance et les efforts de mise en œuvre de la loi peuvent avoir du mal à suivre le rythme. Les décideurs politiques sont confrontés au défi de créer des environnements de gouvernance et/ou de réglementation suffisamment flexibles, agiles et à l'épreuve du temps.
- B. Les développeurs d'IA à usage général sont confrontés à une forte pression concurrentielle, ce qui peut les inciter à mettre en œuvre des mesures d'atténuation des risques moins poussées. Les marchés caractérisés par des coûts fixes élevés, des coûts marginaux faibles et des effets de réseau ont tendance à créer des pressions concurrentielles qui découragent les investissements en matière de sécurité. Le marché de l'IA à usage général est un tel marché.
- C. La croissance rapide et la consolidation du secteur de l'IA suscitent des inquiétudes concernant certaines technologies de l'IA. Les entreprises deviennent particulièrement puissantes parce que des secteurs vitaux de la société dépendent de leurs produits. Ces entreprises peuvent être plus enclines à prendre des risques excessifs ou à négliger les normes de sécurité si elles estiment que la faillite de l'entreprise pourrait coûter cher aux gouvernements.
- D. Le manque inhérent de transparence algorithmique et de transparence institutionnelle
- L'IA à usage général rend la responsabilité juridique difficile à déterminer, ce qui peut entraver la gouvernance et l'application des règles. Le fait que les systèmes d'IA à usage général puissent agir d'une manière qui n'a pas été explicitement programmée ou prévue par leurs développeurs ou utilisateurs soulève des questions quant à savoir qui devrait être tenu responsable des dommages qui en résultent.

Définitions clés

- **Transparence algorithmique** : degré auquel les facteurs qui influencent les résultats de l'IA à usage général, par exemple les recommandations ou les décisions, sont connus par les différentes parties prenantes. Ces facteurs peuvent inclure le fonctionnement interne du modèle d'IA, la manière dont il a été formé, les données sur lesquelles il est formé, les caractéristiques des entrées qui ont affecté ses résultats et les décisions qu'il aurait prises dans des circonstances différentes.
- **Transparence institutionnelle** : la mesure dans laquelle les entreprises d'IA divulguent leurs informations techniques ou des informations organisationnelles soumises à un contrôle public ou gouvernemental, y compris des données de formation, des architectures de modèles, des données sur les émissions, des mesures de sécurité et de sûreté ou des processus décisionnels.
- **Le gagnant rafle tout** : Concept économique faisant référence aux cas dans lesquels une seule entreprise capte une part de marché très importante, même si les consommateurs ne préfèrent que légèrement ses produits ou services à ceux de ses concurrents.
- **Course vers le bas** : un scénario concurrentiel dans lequel des acteurs comme les entreprises ou les États-nations privilégient le développement rapide de l'IA plutôt que la sécurité.

- Avantage du premier entrant : l'avantage concurrentiel obtenu en étant le premier à établir une position de marché significative dans un secteur.
- Formation distribuée : processus de formation de modèles d'IA sur plusieurs processeurs et serveurs, concentrés dans un ou plusieurs centres de données.
- L'humain dans la boucle : une exigence selon laquelle les humains doivent superviser et approuver les processus automatisés dans les domaines critiques.
- Comportement émergent : capacité des systèmes d'IA à agir d'une manière qui n'a pas été explicitement programmée ou prévue par leurs développeurs ou utilisateurs.

A. Alors que les marchés de l'IA à usage général évoluent rapidement, les efforts de gouvernance, de réglementation ou d'application de la loi peuvent avoir du mal à suivre le rythme. Un thème récurrent dans le discours sur le risque de l'IA à usage général est l'inadéquation entre le rythme de l'innovation technologique et le développement des structures de gouvernance (1088). Bien que les cadres juridiques et de gouvernance existants s'appliquent à certaines utilisations des systèmes d'IA à usage général, et que plusieurs juridictions (comme l'Union européenne, la Chine, les États-Unis et le Canada) aient lancé ou achevé des efforts pour établir des normes pertinentes ou pour réglementer l'IA en général et l'IA à usage général en particulier, des domaines d'incertitude réglementaire persistent, en particulier en ce qui concerne les nouvelles capacités de l'IA. Sur un marché qui évolue aussi rapidement que celui de l'IA à usage général, il est très difficile de combler ces lacunes de manière réactive, car au moment où une solution de gouvernance et/ou de réglementation est mise en œuvre, elle peut déjà être obsolète. Par exemple, les critiques de la réglementation des médias sociaux soulignent souvent les défis dans des domaines tels que la confidentialité des données, suggérant que ces problèmes se sont développés plus rapidement que les décideurs politiques n'ont pu les résoudre efficacement (1089, 1090).

Les décideurs politiques sont confrontés au défi de créer des environnements réglementaires flexibles et résistants aux changements technologiques au fil du temps.

Le rythme et l'imprévisibilité des avancées de l'IA à usage général posent un « dilemme de preuves » aux décideurs politiques. Compte tenu des avancées parfois rapides et inattendues, les décideurs politiques devront souvent évaluer les avantages et les risques potentiels des avancées imminentes de l'IA sans disposer d'un large corpus de preuves scientifiques. Ce faisant, ils sont confrontés à un dilemme. D'une part, les mesures préventives d'atténuation des risques fondées sur des preuves limitées pourraient s'avérer inefficaces ou inutiles. D'autre part, l'attente de preuves plus solides d'un risque imminent pourrait laisser la société démunie ou même rendre l'atténuation impossible, par exemple si des progrès soudains dans les capacités de l'IA et les risques associés se produisent. Les entreprises et les gouvernements développent des systèmes d'alerte précoce et des cadres de gestion des risques qui peuvent réduire ce dilemme. Certains d'entre eux déclenchent des mesures d'atténuation spécifiques lorsqu'il existe de nouvelles preuves de risques, tandis que d'autres exigent des développeurs qu'ils fournissent des preuves de sécurité avant de lancer un nouveau modèle.

B. Les développeurs d'IA à usage général sont confrontés à une forte pression concurrentielle, ce qui peut les inciter à mettre en œuvre des mesures d'atténuation des risques moins poussées. Le coût ponctuel du développement d'un modèle d'IA à usage général de pointe est très élevé, tandis que les coûts marginaux de distribution d'un tel modèle à des utilisateurs (supplémentaires) sont relativement faibles. Par exemple, le coût estimé de la formation de GPT-4 était de 40 millions de dollars (27), mais une fois formé, le coût d'exécution du modèle pour une seule requête ne serait que de quelques centimes, ce qui lui permettrait de servir de nombreux utilisateurs à un coût marginal relativement faible. En théorie économique,

Ces conditions peuvent conduire à une dynamique de « gagnant rafle tout » dans laquelle les leaders du secteur peuvent rapidement conquérir un large marché, tandis que les acteurs de deuxième place seront fortement désavantagés. Ainsi, si le fait de prendre des raccourcis (par exemple) dans les tests et la sécurité pouvait permettre à un développeur de prendre la tête en matière de capacité du modèle, il existe alors une forte incitation à prendre des raccourcis (1091). Cette dynamique est visible sur les plateformes de médias sociaux, où une large base d'utilisateurs initiale a attiré plus de personnes à rejoindre certaines plateformes parce que c'est là que se trouvaient leurs amis, ce qui a rendu la plateforme leader plus précieuse pour les nouveaux utilisateurs et a encore élargi son réseau, tandis que les nouveaux réseaux sociaux ont souvent du mal à atteindre une masse critique (1092). La dynamique du « gagnant rafle tout » suscite des inquiétudes quant aux scénarios potentiels de « course vers le bas », où les acteurs rivalisent pour développer des modèles d'IA à usage général le plus rapidement possible tout en sous-investissant dans des mesures visant à garantir que les modèles sont sûrs et éthiques (1093, 1094).

Les marchés caractérisés par des coûts fixes élevés, des coûts marginaux faibles et des effets de réseau ont tendance à créer des pressions concurrentielles qui découragent les investissements en matière de sécurité. La théorie économique et les études empiriques ont montré que, dans des conditions de coûts fixes élevés, de coûts marginaux faibles et d'effets de réseau importants, les entreprises des marchés hautement concurrentiels ont tendance à sous-investir dans les mesures de sécurité (1095, 1096, 1097, 1098). Par exemple, dans les débuts de l'industrie de l'aviation commerciale, les compagnies aériennes opérant avec de faibles marges bénéficiaires en raison des coûts fixes élevés d'acquisition et de maintenance des avions lésinaient parfois sur les procédures de sécurité pour réduire les coûts et maintenir des prix de billets compétitifs (1099). Ces conditions sont présentes sur le marché de l'IA à usage général. De plus, sur les marchés hautement concurrentiels avec des avantages significatifs pour les premiers entrants, la théorie économique suggère que le comportement de prise de risque a tendance à être récompensé et peut devenir répandu parmi les entreprises survivantes (1100). Bien que les études directes sur l'investissement en matière de sécurité sur le marché de l'IA fassent actuellement défaut, ces principes économiques et ces études empiriques dans d'autres domaines suggèrent des motifs d'inquiétude. Cela pourrait contribuer à créer des situations dans lesquelles il est difficile pour les développeurs d'IA à usage général de s'engager unilatéralement à respecter des normes de sécurité strictes, car cela pourrait les placer dans une situation de désavantage concurrentiel (1101). Dans le même temps, d'un point de vue commercial à long terme, la publication de modèles risqués sans mesures de sécurité adéquates pourrait nuire à la confiance des utilisateurs et à la réputation de l'entreprise, créant potentiellement des incitations à investir dans la sécurité plus fortes que ne le suggèrent les pressions concurrentielles à court terme.

C. La croissance et la consolidation rapides du secteur de l'IA suscitent des inquiétudes quant au fait que certaines entreprises d'IA deviennent particulièrement puissantes parce que des secteurs critiques de la société dépendent de leurs produits, ce qui pourrait les inciter à prendre des risques excessifs (voir [2.3.3. Concentration du marché et points de défaillance uniques](#)). De tels scénarios sont bien étudiés dans la littérature économique (1102). Ils surviennent lorsqu'une organisation atteint une taille et un niveau d'influence si importants qu'une défaillance potentielle pourrait présenter des risques systémiques pour l'économie ou la sécurité nationale. Les gouvernements sont donc enclins à prendre des mesures pour protéger ces organisations de la faillite, par exemple en annulant les dettes ou en fournissant des fonds de sauvetage. Lorsqu'elles sont protégées de cette manière, les entreprises peuvent être plus enclines à prendre des risques excessifs ou à faire des concessions sur les normes de sécurité (1103, 1104), bien que les preuves empiriques sur cet effet restent mitigées (1105). Certains craignent que des secteurs critiques de la société ne deviennent au fil du temps trop dépendants des produits d'un petit nombre d'entreprises d'IA de premier plan. Les applications d'IA font de plus en plus partie intégrante de la vie quotidienne et les petites entreprises sont de plus en plus petites.

Les startups cherchent souvent à être rachetées ou à collaborer avec des entreprises de plus grande taille pour surmonter les obstacles à l'entrée sur le marché, notamment les coûts extrêmement élevés de formation d'un modèle d'IA à usage général. Dans de tels accords, les startups échangent généralement l'accès à leurs innovations contre l'utilisation de l'infrastructure informatique et des derniers modèles des grandes entreprises, renforçant encore la concentration du marché et, potentiellement, la dépendance excessive à l'égard des produits d'IA de quelques leaders du secteur (767).

Outre la dynamique de concentration du marché, plusieurs autres facteurs peuvent contribuer au sous-investissement dans l'atténuation des risques. À l'instar de la pollution environnementale ou des problèmes de santé publique tels que le tabac, de nombreux dommages potentiels causés par les systèmes d'IA représentent des externalités – des coûts qui peuvent être supportés par la société plutôt que directement par les développeurs (1106, 1107, 1108). En outre, la théorie économique suggère que lorsqu'il existe un décalage temporel important entre les actions et les conséquences, les acteurs du marché peuvent systématiquement sous-investir dans l'atténuation des risques (1109). Ce défi est aggravé par l'incertitude inhérente à ces dommages potentiels, ce qui rend difficile la quantification du niveau approprié d'investissement dans l'atténuation des risques. Bien que les preuves empiriques sur cette question soient rares, la théorie économique suggère que les coûts immédiats de l'atténuation des risques mis en balance avec les avantages futurs incertains créent des incitations au sous-investissement dans les mesures de sécurité.

D. Le manque de transparence inhérent aux systèmes d'IA à usage général et la transparence institutionnelle limitée des organisations qui développent l'IA rendent la responsabilité juridique difficile à déterminer, ce qui peut entraver la gouvernance et l'application des règles. Le suivi du développement et de l'utilisation des systèmes d'IA est important pour établir la responsabilité des préjudices potentiels, surveiller et rechercher des preuves d'utilisation malveillante et détecter les dysfonctionnements (1002, 1110, 1111). En principe, ce sont les personnes et les entreprises qui sont tenues responsables, et non la technologie, c'est pourquoi les développeurs maintiennent une politique de « l'intervention humaine dans la boucle » pour de nombreux domaines critiques, où un humain doit superviser et approuver des processus par ailleurs automatisés. Cependant, remonter jusqu'aux personnes responsables des préjudices est très difficile (1112, 1113, 1114), tout comme rassembler des preuves d'erreur ou de négligence. Cela découle de facteurs à la fois techniques et institutionnels : les processus décisionnels des modèles d'IA sont difficiles à interpréter, même pour leurs développeurs (manque de transparence algorithmique), et les entreprises d'IA traitent souvent leurs données de formation, leurs méthodologies et leurs procédures opérationnelles comme des informations commercialement sensibles qui ne sont pas ouvertes à l'examen public (manque de transparence institutionnelle) (1025*, 1115, 1116, 1117, 1118, 1119, 1120). Sans transparence des systèmes techniques et des processus organisationnels, il est difficile d'élaborer les types de normes de gouvernance de la sécurité complètes qui sont courantes dans d'autres domaines critiques pour la sécurité tels que l'automobile, les produits pharmaceutiques et l'énergie (1121, 1122, 1123). Le fait que les systèmes d'IA à usage général puissent agir d'une manière qui n'a pas été explicitement programmée ou prévue par leurs développeurs ou utilisateurs soulève des questions quant à savoir qui doit être tenu responsable des dommages qui en résultent (174, 1124). Ces défis en matière de responsabilité deviennent encore plus prononcés avec des systèmes d'IA de plus en plus autonomes qui nécessitent moins de surveillance humaine directe, car il devient plus difficile de relier des actions nuisibles spécifiques aux instructions ou aux décisions humaines ([voir 3.1. Aperçu de la gestion des risques](#)).

La concentration de l'expertise en IA dans les entreprises privées peut créer d'importantes lacunes d'information pour les décideurs politiques et le public. Alors que les chercheurs universitaires et les experts du secteur public contribuent au développement de l'IA et à la recherche sur la sécurité, une grande partie du travail de pointe dans le développement de l'IA se fait au sein des entreprises privées (1125, 1126). Cette concentration d'expertise peut rendre difficile l'accès à l'IA.

Les décideurs politiques et le public ont accès aux connaissances techniques nécessaires pour prendre des décisions éclairées sur la gouvernance de l'IA et la gestion des risques. L'asymétrie d'information qui en résulte entre les développeurs d'IA et les autres parties prenantes pourrait compliquer les efforts visant à élaborer des cadres de gouvernance et/ou de réglementation et des normes de sécurité appropriés.

3.3. Identification et évaluation des risques

INFORMATIONS CLÉS

- L'évaluation des dangers des systèmes d'IA à usage général fait partie intégrante de la gestion des risques.
Les scientifiques utilisent diverses techniques pour étudier les dangers pendant le développement du système, avant le déploiement et après le déploiement.
- Les réglementations et engagements existants en matière d'IA nécessitent une identification et une évaluation rigoureuses des risques.
Les gouvernements et les développeurs d'IA à usage général ont adopté des politiques qui les obligent à identifier et à évaluer les risques et les impacts potentiels des systèmes d'IA à usage général sur les personnes, les organisations et la société.
- Bien que très utiles, les méthodes quantitatives existantes pour évaluer les risques liés à l'IA à usage général ont des limites importantes. Les risques de sécurité dépendent fortement de la manière et du lieu d'utilisation de ces systèmes, qui sont souvent imprévus, ce qui rend difficile la mesure des risques sans deviner comment les gens les utiliseront. Cela est particulièrement difficile pour l'IA à usage général, car elle peut être utilisée dans d'innombrables situations différentes et de nombreux préjudices potentiels (par exemple, biais, toxicité et désinformation) sont difficiles à mesurer objectivement. Bien que les méthodes actuelles d'évaluation des risques en soient à leurs balbutiements, elles peuvent être grandement améliorées.
- Une évaluation rigoureuse des risques nécessite de combiner plusieurs approches d'évaluation,
Les principaux indicateurs de risque incluent les évaluations des systèmes eux-mêmes, la manière dont les personnes les appliquent, ainsi que l'analyse prospective des menaces. Pour que les évaluations à la frontière technique soient efficaces, les évaluateurs doivent disposer de capacités et d'une expertise techniques substantielles et croissantes. Ils ont également besoin de suffisamment de temps et d'un accès plus direct que celui actuellement disponible aux modèles, aux données de formation, aux méthodologies utilisées et aux évaluations internes à l'entreprise – mais les entreprises qui développent une IA à usage général ne sont généralement pas fortement incitées à les accorder.
- Ces derniers mois, de plus en plus de recherches ont évalué l'efficacité des méthodes d'évaluation des risques par l'IA, en identifiant les lacunes actuelles et les critères d'amélioration. Bien que davantage de données soient nécessaires – en particulier pour les nouveaux risques – ces progrès techniques sont complétés par des développements institutionnels, à mesure que les gouvernements commencent à renforcer leurs capacités d'évaluation et que les parties prenantes s'efforcent d'établir des lignes directrices plus claires sur les personnes responsables des différents aspects de l'évaluation des risques.
- L'absence de normes claires d'évaluation des risques et d'évaluations rigoureuses crée un défi politique urgent, car les modèles d'IA sont déployés plus rapidement que les risques qu'ils présentent ne peuvent être évalués. Les décideurs politiques sont confrontés à deux défis majeurs : 1. les évaluations internes des risques par les entreprises sont essentielles pour la sécurité mais insuffisantes pour une surveillance adéquate, et 2. les audits complémentaires effectués par des tiers et des organismes réglementaires nécessitent davantage de ressources, d'expertise et d'accès aux systèmes que ce qui est actuellement disponible.

Définitions clés

- **Risque** : la combinaison de la probabilité et de la gravité d'un préjudice résultant du développement, du déploiement ou de l'utilisation de l'IA.
- **Danger** : tout événement ou activité susceptible de causer un préjudice, tel qu'une perte de vie, une blessure, une perturbation sociale ou des dommages environnementaux.
- **Déploiement** : processus de mise en œuvre de systèmes d'IA dans des applications, des produits ou des services du monde réel où ils peuvent répondre aux demandes et fonctionner dans un contexte plus large.
- **Évaluations** : Évaluations systématiques des performances, des capacités et des performances d'un système d'IA. vulnérabilités ou impacts potentiels. Les évaluations peuvent inclure des analyses comparatives, des équipes rouges et des audits et peuvent être menées avant et après le déploiement du modèle.
- **Benchmark** : un test ou une mesure standardisé, souvent quantitatif, utilisé pour évaluer et comparer les performances des systèmes d'IA sur un ensemble fixe de tâches conçues pour représenter le monde réel usage.
- **Red-teaming** : Un processus systématique dans lequel des individus ou des équipes dédiés recherchent vulnérabilités, limitations ou potentiel d'utilisation abusive par diverses méthodes. Souvent, l'équipe rouge recherche des entrées qui induisent un comportement indésirable dans un modèle ou un système pour identifier les failles de sécurité.
- **Jailbreaking** : Génération et envoi d'invites conçues pour contourner les garde-fous et amener un système d'IA à produire du contenu nuisible, tel que des instructions pour la construction d'armes.
- **Audit** : Un examen formel de la conformité d'une organisation aux normes, aux politiques et procédures, généralement effectuées par un tiers indépendant.
- **Rapports d'incidents** : documenter et partager les cas dans lesquels le développement ou le déploiement de l'IA a causé des préjudices directs ou indirects.

Pour gérer les risques liés à l'IA à usage général, il est nécessaire de comprendre et de mesurer les risques qu'elle représente pour les personnes, les organisations et la société. Plusieurs gouvernements et développeurs d'IA à usage général ont déjà adopté des politiques et des réglementations qui les obligent à identifier et à évaluer les risques et les impacts potentiels des systèmes d'IA à usage général, déclenchant des réponses planifiées lorsque les risques atteignent des seuils spécifiques. L'« identification des risques » est le processus d'identification des risques potentiels de la technologie, y compris les dangers possibles et les résultats imprévus. L'« évaluation des risques » est le processus d'évaluation de la gravité et de la probabilité d'occurrence de chaque risque identifié. (Voir le tableau 3.1 dans [3.1 Présentation de la gestion des risques](#) pour un aperçu des étapes de gestion des risques, y compris [l'identification et l'évaluation des risques](#) ainsi que l'évaluation des risques, l'atténuation des risques et la gouvernance des risques).

Méthodes d'identification des risques

Les risques généraux liés à l'IA peuvent être identifiés et formulés à différents niveaux de spécificité. Par exemple, une grande catégorie de risques généraux liés à l'IA est la fabrication ou les « hallucinations ». désinformation – c'est-à-dire la génération de résultats inexacts ou trompeurs. Un exemple plus spécifique du même risque est celui d'une IA polyvalente qui invente un lieu de vote inexistant lorsque l'utilisateur lui demande de recueillir des informations sur l'endroit où voter lors d'une élection nationale.

(1127). La spécification d'un risque peut faciliter ou compliquer la tâche des évaluateurs en ce qui concerne sa gravité et sa probabilité. Les risques mieux spécifiés sont plus faciles à évaluer et à atténuer.

Les évaluateurs doivent bien comprendre les cas d'utilisation de l'IA à usage général afin de conceptualiser ses risques avec le degré de spécificité approprié. Par exemple, si les utilisateurs d'IA à usage général sont susceptibles de l'inciter à recueillir des informations sur les campagnes politiques et les procédures de vote, l'évaluation du risque que le modèle « hallucine un lieu de vote » peut être une priorité élevée. Par conséquent, les approches participatives, qui consistent à s'engager auprès de diverses parties prenantes et communautés concernées pour comprendre leurs cas d'utilisation, leurs pratiques, leurs besoins et leurs valeurs, sont particulièrement utiles pour identifier les risques prioritaires pour les utilisateurs. Les audits de foule (1128) sont un exemple d'approche participative. Ils sont conçus pour permettre aux utilisateurs quotidiens de faire émerger de manière collaborative les dommages potentiels des produits et services d'IA. La création de mécanismes accessibles permettant au public de signaler les dommages observés et perçus est une autre méthode importante d'identification des risques. Les bases de données de suivi des incidents d'IA, telles que l'AI Incidents Monitor (AIM) de l'OCDE, sont des plateformes conçues pour collecter, catégoriser et signaler les incidents préjudiciables impliquant l'IA (459). En bref, il est nécessaire d'identifier et d'évaluer les risques dans ce contexte.

Pour faciliter les pratiques d'identification des risques de l'IA à usage général, les chercheurs ont proposé des taxonomies de dangers (439, 933, 951*, 1129). Ces taxonomies répertorient les catégories de risques, telles que les dangers informationnels, la mémorisation des données de formation (qui peut conduire à une violation du droit d'auteur et à des problèmes de confidentialité) et l'utilisation malveillante (par exemple, l'écriture de logiciels malveillants). Les taxonomies de dangers peuvent servir de point de départ pour aider les évaluateurs à conceptualiser, identifier et spécifier les risques saillants associés à l'IA à usage général dans des domaines d'application spécifiques. Dans la gestion des risques et l'ingénierie de la sécurité conventionnelles, il existe plusieurs méthodes bien établies pour identifier les dangers et les risques d'une technologie, notamment l'analyse des défaillances fonctionnelles et l'HAZOP (étude des dangers et de l'opérabilité) (1130). Ces méthodes ont été adoptées dans un large éventail d'industries, y compris l'industrie automobile, qui prend également en compte la SOTIF (946). Outre les typologies et taxonomies des risques, des travaux récents ont commencé à adapter certaines de ces techniques conventionnelles, par exemple l'analyse des risques, la méthode du nœud papillon et les études de sécurité, aux produits et services d'IA (968, 1131, 1132, 1133), mais des recherches supplémentaires sont nécessaires dans ce domaine. Voir [3.1 Aperçu de la gestion des risques](#) pour une discussion sur d'autres pratiques d'identification des risques établies dans d'autres domaines.

Méthodes d'évaluation des risques

Une fois les risques hautement prioritaires identifiés, ils doivent être évalués pour déterminer la probabilité et la gravité du préjudice, du danger ou du résultat imprévu en question.

Il est essentiel de mieux comprendre l'état actuel des méthodes d'évaluation des risques de l'IA à usage général pour la politique en la matière, car les évaluations des risques sont un élément essentiel de nombreuses approches de gouvernance et de réglementation de l'IA. Par exemple, la loi européenne sur l'IA classe les systèmes d'IA en quatre niveaux de risque principaux en fonction de leur impact potentiel et impose des exigences différentes aux systèmes d'IA en fonction de leur niveau de risque. En outre, de nombreuses entreprises leaders dans le domaine de l'IA ont accepté de créer des engagements en matière de sécurité de l'IA avec

Les mesures d'atténuation doivent être proportionnelles et spécifiques au risque évalué de leurs systèmes (1134). Cependant, l'évaluation des risques est un sujet de recherche relativement nouveau dans la communauté de la sécurité de l'IA, et il n'existe actuellement aucune approche systématique entièrement validée pour évaluer la gravité et la probabilité des dommages causés par l'IA à usage général. La mise en œuvre des politiques susmentionnées nécessitera un domaine d'évaluation des risques beaucoup plus mature pour l'IA à usage général.

Les travaux existants sur la sécurité de l'IA se concentrent principalement sur les approches conventionnelles de test de modèles en IA, souvent menées après le développement de modèles d'IA à usage général. Ce recours à une évaluation rétrospective (par opposition à prospective) des risques peut conduire à des omissions majeures et à des estimations erronées des risques hautement prioritaires. Dans la gestion des risques et l'ingénierie de la sécurité conventionnelles, une étape critique de l'évaluation des risques est l'analyse prospective des risques avant d'achever la conception et le développement d'un système. Cette étape est actuellement souvent négligée dans les évaluations des risques de l'IA à usage général. Dans la sécurité de l'IA, l'évaluation des risques consiste principalement à exécuter une batterie de tests et d'évaluations sur le système d'IA à usage général, puis à traduire les résultats en estimations quantitatives des risques. Cela contraste avec l'évaluation des risques traditionnelle, qui consiste à 1. analyser les causes, les conséquences et la prévalence des risques (par des méthodes telles que la cartographie causale et la technique Delphi) puis 2. évaluer si le risque est acceptable, par exemple au moyen de listes de contrôle et de matrices de risques. Des travaux récents ont commencé à adapter certaines de ces techniques aux produits et services d'IA (944, 968). Voir [3.1. Aperçu de la gestion des risques](#) pour une discussion plus approfondie des approches d'évaluation des risques établies dans d'autres domaines.

Les approches et méthodologies techniques existantes pour l'évaluation des risques de l'IA à usage général s'appuient largement sur des tests et des évaluations qui peuvent être divisés en quatre couches (1135) :

1. Les tests de modèle évaluent le modèle d'IA à usage général en termes de (souvent quantitatifs)
Mesures de performance sur des tâches proxy conçues pour représenter une utilisation dans le monde réel. Ces tests prennent souvent la forme de tests de performance, c'est-à-dire d'ensembles fixes d'invites sur lesquelles tester un modèle.
2. Le red teaming est un processus systématique dans lequel des individus ou des équipes dédiées recherchent des vulnérabilités, des limites ou un potentiel d'utilisation abusive dans les modèles ou les systèmes d'IA par le biais de diverses méthodes.
Souvent, l'équipe rouge recherche des entrées qui induisent un comportement indésirable dans le but d'améliorer les protections du modèle ou du système contre de telles attaques.
3. Les tests sur le terrain évaluent les risques de l'IA à usage général dans des conditions réelles.
4. Les évaluations d'impact à long terme surveillent et évaluent les impacts à long terme du système sur les personnes, les organisations et la société.

L'une des principales lacunes en matière de données probantes est la recherche visant à établir la validité, la fiabilité et la praticabilité des méthodes d'évaluation des risques de l'IA à usage général existantes. Les bonnes méthodes de mesure des risques doivent être valides, fiables et pratiques. La validité fait référence à la mesure dans laquelle un test, un outil ou un instrument mesure avec précision ce qu'il est censé mesurer. Par exemple, des problèmes de validité surviennent si une référence diffère de l'utilisation dans le monde réel ou contient de fausses étiquettes (1136). La fiabilité fait référence à la cohérence, à la stabilité et à la fiabilité d'une mesure au fil du temps et dans différents contextes. En d'autres termes, elle indique le degré auquel une mesure produit des résultats cohérents et fiables.

résultats reproductibles dans des conditions similaires (1137). Des travaux antérieurs ont montré que même de petites perturbations des invites peuvent avoir des effets significatifs sur le comportement et les performances de l'IA à usage général sur les benchmarks (1138, 1139). La praticité évalue si la mesure peut être effectuée de manière efficace et efficiente dans la pratique par les évaluateurs désignés, en tenant compte de contraintes telles que le temps, le coût, la disponibilité des ressources informatiques et la charge de travail des évaluateurs. Par exemple, le processus d'évaluation de l'IA à usage général repose de plus en plus sur l'utilisation de l'IA à usage général (522, 929*), ce qui nécessite des capacités techniques et soulève de nouvelles préoccupations (par exemple concernant les agents LLM privilégiant les résultats de leur propre famille de modèles (1140). Pour une évaluation rigoureuse des risques, la validité et la fiabilité sont prioritaires sur la facilité et la commodité de la mesure (1141).

Depuis la publication du rapport intermédiaire, la communauté scientifique a progressé dans la mise en œuvre et l'évaluation des méthodes d'évaluation des risques existantes. Les instituts de sécurité de l'IA des États-Unis et du Royaume-Uni (US AISI et UK AISI) ont récemment publié un rapport technique détaillant une évaluation préalable au déploiement de la version améliorée de Claude 3.5 Sonnet (1142). De nouvelles recherches ont examiné la reproductibilité (1143, 1144*) ou la validité, qui peuvent être compromises lorsque les modèles d'IA sont formés ou exposés à des données de test au préalable (contamination de référence) (1145, 1146). Cependant, des preuves supplémentaires sont nécessaires pour caractériser les forces et les faiblesses des méthodes d'évaluation de l'IA à usage général existantes (465), en particulier lorsque l'IA à usage général est utilisée dans de nouveaux domaines.

La couche initiale de l'évaluation des risques de l'IA à usage général consiste souvent à tester le comportement du modèle sur certaines tâches de référence fixes. De nouveaux tests de référence et des tests standardisés

Des mesures ont été conçues pour évaluer et comparer diverses catégories de risques pour les applications d'IA à usage général dans des scénarios et des tâches stylisés (122, 137, 141, 1147*, 1148, 1149*). Par exemple, l'AI Safety Benchmark de MLCommons (457) fournit une référence pour mesurer sept catégories de risques, telles que la désinformation et le contenu préjudiciable. L'évaluation holistique des modèles linguistiques (HELM) comprend 16 scénarios et sept mesures, dont la robustesse, l'équité et le biais (1150). Les évaluations des capacités nuisibles (318*) sont utilisées pour évaluer si l'IA à usage général possède des connaissances ou des compétences particulièrement dangereuses (telles que la capacité à aider aux cyberattaques (2.1.3. [Cyberinfraction](#)) ou à aider à la conception d'armes biologiques (2.1.4. [Attaques biologiques et chimiques](#))). Les décisions importantes à venir des entreprises et des gouvernements concernant la diffusion des modèles reposent en partie sur ces évaluations (596*, 947*, 1134). Les référentiels existants varient considérablement en qualité (1151) et leur champ d'application est souvent flou. Certaines bonnes pratiques pour créer des référentiels de haute qualité ont été proposées (1151, 1152*).

Si les méthodes de test de modèles peuvent constituer une première étape nécessaire pour évaluer les risques de l'IA à usage général, elles ne suffisent pas à elles seules. Il est impossible de tirer des conclusions quantitatives fiables sur les risques que ces méthodes visent à saisir sans formuler de solides hypothèses sur les modèles d'utilisation dans des applications spécifiques. De telles hypothèses sont difficiles à justifier : tout d'abord, la technologie est à usage général et peut être utilisée dans de nombreux contextes, il est donc difficile de prédire les modèles d'utilisation. Ensuite, certains risques (par exemple, les biais, la toxicité et la désinformation) sont difficiles à évaluer.

Les critères de référence ne peuvent pas être définis de manière objective, et toute définition doit reposer sur des hypothèses discutables sur ce qui est (par exemple) « toxique » ou « biaisé ». Par conséquent, les critères de référence ne peuvent pas saisir les risques associés à l'utilisation de l'IA à usage général dans de nouveaux domaines et pour de nouvelles tâches, car les conditions de test diffèrent toujours de l'utilisation dans le monde réel à des degrés divers (1153*). Les critères de référence servent au mieux de mesure indirecte de la catégorie de risque en question (par exemple, les évaluations subjectives des annotateurs humains ou des modérateurs de contenu peuvent servir de mesure indirecte de la « toxicité » (1154)). Cependant, ces mesures indirectes ne reflètent souvent pas de manière fiable le risque réel dans le contexte. Par exemple, si les évaluateurs humains ne sont pas diversifiés, cela peut conduire à des critères de référence contenant des étiquettes biaisées, car des personnes issues d'horizons similaires peuvent systématiquement manquer certains exemples de toxicité ou de désinformation. De plus, l'amélioration des scores d'un critère de référence ne se traduit pas toujours par une diminution du risque associé dans la pratique. Par exemple, un LLM peut réussir l'examen du barreau pour les avocats, mais cela ne signifie pas qu'il peut créer des notes juridiques efficaces (445, 446, 451). Tout référentiel fixe est souvent facile à améliorer sans atténuer le risque cible (1070). Si la création de capacités pour des référentiels collaboratifs en évolution dynamique peut répondre à certains de ces défis, il est important que les évaluateurs de l'IA comprennent les limites inhérentes aux approches quantitatives des tests de modèles (1155) et évitent de s'y fier excessivement comme couche principale d'évaluation des risques.

Le red-teaming et les attaques adverses sont d'autres méthodes courantes pour identifier et évaluer les risques, mais peuvent nécessiter un accès spécial. Le terme « red-team » désigne un ensemble d'évaluateurs chargés de trouver les vulnérabilités d'un système en l'attaquant. Contrairement aux tests de performance, qui sont pour la plupart statiques et consistent en un ensemble fixe de cas de test, l'un des principaux avantages du red-teaming est qu'il adapte l'évaluation au système spécifique testé. Grâce aux interactions adverses avec un système, les red-teamers peuvent concevoir des entrées personnalisées pour identifier les comportements les plus défavorables, les opportunités d'utilisation malveillante et les pannes inattendues. À titre d'exemple, les attaques contre les modèles de langage peuvent prendre la forme d'entrées générées automatiquement (904*, 1053, 1063*, 1156, 1157, 1158*, 1159, 1160, 1161, 1162) ou générées manuellement (1056*, 1059, 1158*, 1163). Dans les attaques automatisées, par exemple, les LLM peuvent être utilisés pour générer des invites conçues pour amener un autre système d'IA à produire du contenu nuisible, comme des instructions sur des matières dangereuses, même après un refus initial du système. Ces attaques de « jailbreaking » contournent les restrictions de sécurité des modèles (460, 904*, 1052, 1053, 1164, 1165*). Les approches automatisées peuvent tester systématiquement des milliers de variantes d'attaques potentielles, ce qui permet une couverture plus étendue et plus rapide que les seuls tests manuels. Cependant, la red-teaming manuelle sur des conversations plus longues peut détecter des problèmes que les attaques automatisées actuelles peuvent manquer (1056*). Cependant, elle peut être lente, demander beaucoup de travail et nécessiter un accès spécial. Des recherches supplémentaires pour une red-teaming automatisée plus rapide et efficace sont nécessaires pour relever ce défi (1166).

Bien que le red-teaming soit plus efficace pour faire apparaître un plus large éventail de risques généraux liés à l'IA que les tests de modèles, de nombreux préjudices et dangers importants peuvent rester non détectés. Il est important de noter que si une activité de red-teaming ne parvient pas à faire apparaître certaines catégories de risques, cela ne signifie pas que ces risques sont peu probables. Des travaux antérieurs ont montré que les bugs échappent souvent à la détection (1022). Un exemple concret est celui des jailbreaks, qui incitent les systèmes de chat à usage général à se conformer aux demandes nuisibles qu'ils ont été formés à refuser (460, 904*, 1052, 1053, 1164), et qui échappent à la détection initiale par les développeurs (48*, 147*, 1158*). La recherche a également remis en question la capacité du red-teaming à détecter les risques.

Les résultats obtenus sont fiables et reproductibles. Une étude montre que les pratiques de red-teaming dans l'industrie divergent selon plusieurs axes clés, notamment le cadre (par exemple, les caractéristiques des red-teamers et les ressources et méthodes dont ils disposent) et les décisions qu'il éclaire (par exemple, les rapports ultérieurs, la divulgation et l'atténuation) (1167). La composition de l'équipe rouge et les instructions fournies aux red-teamers (1168*), le nombre de rounds d'attaque (1056*) et la disponibilité d'outils auxiliaires ou d'automatisation (1161, 1169) peuvent influencer considérablement les résultats de l'activité, y compris la surface de risque couverte. Voir le tableau 3.2 pour un aperçu des critères de structuration des activités de red-teaming dans la pratique. Des lignes directrices complètes sur le red-teaming visent à relever certains de ces défis (1170).

Phase	Questions et considérations clés
0. Critères de pré-activité	Quel est l'artefact évalué via l'activité de red-teaming proposée ?
	Quel est le modèle de menace que l'activité de red-teaming vise à recréer ?
	Quelle est la vulnérabilité spécifique que l'activité de red-teaming vise à trouver ?
	Quels sont les critères permettant d'évaluer le succès de l'activité de red-teaming ?
	Quelle est la composition de l'équipe, ou qui fera partie de l'équipe ?
1. Critères intra-activité	Quelles ressources sont à la disposition des participants ?
	Quelles instructions sont données aux participants pour guider l'activité ?
	Quel type d'accès les participants ont-ils au modèle ?
	Quelles méthodes les membres de l'équipe peuvent-ils utiliser pour tester l'artefact ?
2. Critères post-activité	Quels rapports et documents sont produits sur les résultats de l'activité ?
	Quelles ressources ont été consommées par l'activité ?
	Dans quelle mesure l'activité a-t-elle été réussie par rapport aux critères spécifiés dans la phase 0 ?
	Quelles sont les mesures proposées pour atténuer les risques identifiés en phase 1 ?

Tableau 3.2 : Différents types de critères peuvent aider les praticiens à structurer le red-teaming avant, pendant et après les activités concernées. Source : basé sur les critères proposés par Feffer et al., 2024 (1167).

Les « tests sur le terrain » sont des exercices conçus pour évaluer les risques dans des conditions d'utilisation normales. Les « études d'amélioration de l'efficacité humaine » examinent si les individus peuvent utiliser l'IA pour effectuer des tâches malveillantes mieux qu'ils ne le pourraient sans l'IA. Les études d'amélioration de l'efficacité humaine sont une variante importante des tests sur le terrain. Elles visent à mesurer comment l'accès à des systèmes d'IA à usage général améliore les compétences et les performances des individus. Par exemple, une étude sur l'amélioration des capacités humaines pourrait explorer la manière dont un système d'IA affecte la capacité d'une personne à accomplir des tâches complexes, telles que le support client (662) ou les opérations de cybersécurité (potentiellement dangereuses) (361, 1171, 1172, 1173), par rapport à ses performances sans l'aide de l'IA. Ces études visent à quantifier l'« amélioration » des capacités humaines et à évaluer si le soutien de l'IA introduit de nouveaux risques, tels que l'abaissement des obstacles à une conduite préjudiciable (voir [2.4. Impact des modèles d'IA polyvalents à pondération ouverte sur les risques liés à l'IA pour une discussion plus approfondie](#) des études sur l'amélioration). Cependant, la conception et la conduite de telles études posent plusieurs défis, notamment la simulation de conditions similaires

Les évaluateurs pourraient relever certains de ces défis s'il existait de meilleures directives pour mener des études d'amélioration de la sécurité sur l'homme et les intégrer dans le déploiement progressif de produits d'IA à usage général. Dans d'autres secteurs critiques pour la sécurité, par exemple les tests de médicaments dans les essais cliniques, une série d'études sont menées dans des conditions de plus en plus réalistes (par exemple, en passant des tests sur les animaux aux études sur les sujets humains), avant que le médicament ne soit jugé prêt à être commercialisé. Une approche similaire pourrait s'avérer utile pour développer des méthodes de test sur le terrain efficaces pour l'IA à usage général.

Français Certains risques associés à l'IA à usage général ne se manifesteront probablement qu'à long terme, ce qui rend les évaluations d'impact à long terme cruciales. Ces risques comprennent les effets de la technologie sur les marchés du travail et l'avenir du travail (2.3.1. [Risques liés au marché du travail](#)), les risques associés à des futurs systèmes d'IA plus performants (2.2.3. [Perte de contrôle](#), 2.1.3. [Cyberinfraction](#), 2.1.4. [Attaques biologiques et chimiques](#)), l'impact environnemental du développement et de l'utilisation de l'IA (voir 2.3.4. [Risques pour l'environnement](#)) et les impacts à long terme sur la cognition, le bien-être et le contrôle humains (1003). Une surveillance, une enquête et une rectification minutieuses des dommages à long terme sont nécessaires pour maintenir la confiance du public dans la technologie et éviter les appels à des contrôles inutilement stricts. Il est difficile d'évaluer avec précision les impacts sociétaux en aval de l'IA à usage général en raison 1. des incertitudes entourant les capacités des futurs systèmes d'IA à usage général, et 2. de l'existence de nombreux facteurs de confusion qui rendent difficile l'attribution des tendances à long terme à une cause unique. La création de capacités permettant de prédire et de surveiller les impacts sociétaux potentiels en aval de l'IA à usage général nécessite une analyse multidisciplinaire et l'implication de perspectives diverses (929*, 1174, 1175).

Défis et opportunités

Outre les défis évoqués ici, voir également 3.2.1. [Défis techniques pour la gestion des risques et l'élaboration des politiques](#) et 3.2.2. [Défis sociétaux pour la gestion des risques et l'élaboration des politiques](#)

La culture du « construire puis tester » en IA entrave l'évaluation et l'atténuation complètes des risques.

Dans la gestion des risques conventionnelle, l'évaluation des risques est intégrée à toutes les étapes de la conception, du développement et du déploiement des produits, et est étroitement liée aux stratégies d'atténuation des risques. En matière de sécurité de l'IA, cependant, les méthodes actuelles d'évaluation des risques sont en grande partie réalisées après le développement et indépendamment de l'atténuation des risques. Des travaux antérieurs (978) ont proposé la création d'études de cas de sécurité et de garanties de sécurité pour l'IA (1176). L'adaptation et la mise en œuvre de telles pratiques pour l'IA à usage général nécessitent à la fois un changement culturel et des recherches plus approfondies.

Les quatre niveaux d'évaluation des risques (tests de modèles, red-teaming, tests sur le terrain et évaluation de l'impact à long terme) sont nécessaires mais pas suffisants pour une évaluation complète des risques. Les méthodes existantes ne fournissent pas de garanties ou d'assurances généralisables concernant la probabilité et la gravité des dommages causés par l'IA à usage général (1177). Les principales lacunes en matière de données probantes concernent 1. l'évaluation de la validité, de la fiabilité et de la praticabilité de chaque niveau d'évaluation de manière indépendante, et 2. la combinaison des informations provenant de différents niveaux d'évaluation pour produire des informations exploitables (41).

En pratique, la réalisation d'une évaluation complète des risques nécessite un accès, des ressources et un temps considérables, qui sont souvent limités. Très peu d'entités disposent des ressources (ou de la volonté d'allouer les ressources nécessaires) pour mener des évaluations complètes, et les conflits d'intérêts potentiels peuvent conduire à des résultats et des rapports trompeurs (1014, 1178). De plus, les évaluateurs ne disposent parfois pas de suffisamment de temps pour tester minutieusement les modèles. Dans certains cas, les entreprises n'ont accordé aux évaluateurs que quelques jours pour tester un nouveau modèle avant sa publication (2*, 129). Une évaluation efficace des modèles nécessite beaucoup de temps et de ressources.

De plus, les développeurs de systèmes d'IA à usage général de pointe limitent souvent l'accès externe à leur technologie (880). Pour les modèles hébergés sur la plateforme d'un développeur ou auxquels il faut accéder via une API (donnant un accès « boîte noire », uniquement aux entrées et sorties du modèle), il peut être difficile pour les évaluateurs externes d'effectuer des attaques adverses, des interprétations de modèles et des réglages fins efficaces (1086, 1179). Par exemple, les modèles d'IA sont généralement formés pour refuser les requêtes dangereuses, mais pour évaluer les capacités dangereuses, les évaluateurs doivent avoir accès aux versions du modèle sans cette barrière de sécurité. Cet accès est parfois fourni (2*). Sans cela, certains risques hautement prioritaires peuvent être négligés. Des informations incomplètes sur la manière dont un système a été conçu, y compris les données, les techniques, les détails de mise en œuvre et les détails organisationnels, entravent les évaluations du processus de développement (34, 488, 1086, 1180, 1181, 1182). Certains chercheurs ont fait valoir qu'une combinaison de mesures techniques, physiques et juridiques peut offrir un accès direct aux chercheurs externes sans compromettre les secrets commerciaux plus qu'ils ne le sont déjà (1086). Plusieurs études ont préconisé des « ports sûrs » juridiques (1036) ou des régimes d'accès médiatisés par le gouvernement (939) pour permettre aux évaluateurs de mener des évaluations indépendantes sans risquer d'être poursuivis ou interdits d'utilisation. Les chercheurs ont proposé des méthodes d'accès structuré qui ne nécessitent pas de rendre publics le code du modèle et les pondérations d'entraînement (1183), mais qui permettent aux chercheurs et auditeurs indépendants d'accéder pleinement au modèle dans un environnement sécurisé conçu pour éviter les fuites. Les chercheurs développent des techniques d'audit qui utilisent des « enclaves sécurisées ». Ces techniques ont le potentiel d'éviter la fuite des paramètres du modèle aux auditeurs, ainsi que des détails de l'audit aux développeurs du modèle (1184).

Une évaluation réussie des risques nécessite la participation de perspectives diverses au processus d'évaluation. La composition de l'équipe d'évaluation en couches d'évaluation, telles que le red-teaming, peut jouer un rôle essentiel dans le processus de découverte, de caractérisation et de priorisation des préjudices (1185). L'amélioration de la participation des parties prenantes a été au centre des préoccupations de la communauté de l'apprentissage automatique ces dernières années (932, 1186, 1187). De nombreuses stratégies ont été proposées, allant de l'élargissement de la compréhension des « impacts » dans les évaluations d'impact de l'IA (1188) à la possibilité d'une gamme plus inclusive de retours d'information humains (1189, 1190). Cependant, pour favoriser la participation, il faut être sensible à plusieurs critères (1186), comme le respect des parties prenantes afin de minimiser le risque d'exploitation (540) et la mise en évidence de choix difficiles entre des valeurs ou des priorités incompatibles (467, 538, 574). Ce processus peut être facilité par des méthodes issues de l'éthique pratique telles que « l'équilibre réflexif » – l'ajustement mutuel des principes et des jugements jusqu'à ce qu'ils soient en accord les uns avec les autres (1191).

Les décideurs politiques sont confrontés à plusieurs défis pour encourager l'identification et l'évaluation adéquates des risques pour les systèmes d'IA à usage général. En l'absence de lignes directrices, de normes et de ressources claires concernant l'évaluation des risques liés à l'IA à usage général, les praticiens sont confrontés à des incertitudes quant à ce qui constitue une évaluation adéquate des risques dans leurs cas d'utilisation spécifiques. Il est donc difficile pour les décideurs politiques d'encourager la conformité. Un autre défi politique est de savoir comment attribuer la responsabilité des différents niveaux d'évaluation des risques aux différents groupes de parties prenantes de l'IA à usage général, notamment les créateurs de technologies, les utilisateurs et les auditeurs tiers (763). Une autre approche consiste à créer des ressources (par exemple, des « bacs à sable » et des « ports sûrs ») qui favorisent les évaluations d'intérêt public (1036) ou les audits par des tiers. Le succès de cette approche dépend fortement de la disponibilité des ressources, d'évaluateurs et d'experts formés, d'incitations à mener des évaluations rigoureuses (par exemple, en offrant une indemnité et une compensation) et d'un accès à des modèles ou à des informations sur les données et les méthodes utilisées. Plusieurs gouvernements ont commencé à renforcer les capacités de réalisation d'évaluations et d'audits techniques de l'IA à usage général. Il reste à voir dans quelle mesure ces efforts feront progresser l'interdisciplinarité et l'évaluation inclusive de l'IA à usage général dans un avenir proche, et dans quelle mesure ils pourront et seront mis à l'échelle dans la pratique (537, 540, 1192, 1193).

3.4. Atténuation et surveillance des risques

3.4.1. Former des modèles plus fiables

INFORMATIONS CLÉS

- Les méthodes de formation actuelles montrent des progrès dans l'atténuation des risques de sécurité liés aux dysfonctionnements et aux utilisations malveillantes, mais restent fondamentalement limitées. Des progrès ont été réalisés dans la formation de modèles d'IA à usage général pour fonctionner de manière plus sûre, mais aucune méthode actuelle ne peut empêcher de manière fiable même les actions manifestement dangereuses.
- Une approche à plusieurs volets apparaît comme nécessaire pour la sécurité.
La fiabilité des modèles nécessite l'analyse de nombreux aspects de leur comportement et de leur processus de développement – notamment l'exactitude des faits, la qualité de la supervision humaine, les mécanismes internes du système d'IA et l'analyse des schémas d'utilisation potentiellement abusive – autant d'éléments qui doivent éclairer les méthodologies de formation. Bien qu'il existe des techniques permettant de supprimer les capacités nuisibles, les méthodes actuelles tendent à les supprimer plutôt qu'à les éliminer.
- L'entraînement antagoniste offre une robustesse limitée contre les attaques. L'entraînement antagoniste consiste à exposer délibérément les modèles d'IA à des exemples conçus pour les faire échouer ou mal se comporter pendant l'entraînement, dans le but de renforcer la résistance à de tels cas. Cependant, les adversaires peuvent toujours trouver de nouveaux moyens (« attaques ») pour contourner ces mesures de protection avec un effort faible à modéré, comme les « jailbreaks » qui conduisent les modèles à se conformer à des demandes nuisibles même s'ils ont été réglés pour ne pas le faire.
- Depuis la publication du rapport intermédiaire (mai 2024), les avancées récentes révèlent à la fois des progrès et de nouvelles préoccupations. Une meilleure compréhension des mécanismes internes des modèles a permis de faire progresser à la fois les attaques et les défenses adverses sans vainqueur clair. De plus, de plus en plus d'éléments suggèrent que les méthodes de formation actuelles – qui reposent largement sur des commentaires humains imparfaits – amènent par inadvertance les modèles à induire les humains en erreur sur des questions difficiles en rendant les erreurs plus difficiles à repérer. L'amélioration de la quantité et de la qualité des commentaires humains est une voie de progrès, bien que les techniques de formation naissantes utilisant l'IA pour détecter les comportements trompeurs soient également prometteuses.
- Les principaux défis auxquels sont confrontés les décideurs politiques concernent l'incertitude et la vérification. Il n'existe pas de méthodes fiables pour quantifier le risque de défaillances inattendues des modèles. Si certains chercheurs explorent des approches dont la sécurité est prouvée, celles-ci restent théoriques. Cela suggère que les cadres de formation à la sécurité doivent actuellement se concentrer sur les processus de recherche, de réponse et d'atténuation des nouvelles défaillances avant qu'elles ne causent des dommages inacceptables.

Définitions clés

- **Interprétabilité** : degré auquel les humains peuvent comprendre le fonctionnement interne d'un modèle d'IA, notamment pourquoi il a généré un résultat ou une décision particulière. Un modèle est hautement interprétable si ses processus mathématiques peuvent être traduits en concepts qui permettent aux humains de retracer les facteurs et la logique spécifiques qui ont influencé le résultat du modèle.
- **Red-teaming** : Un processus systématique dans lequel des individus ou des équipes dédiés recherchent des vulnérabilités, limitations ou potentiel d'utilisation abusive par diverses méthodes. Souvent, l'équipe rouge recherche des entrées qui induisent un comportement indésirable dans un modèle ou un système pour identifier les failles de sécurité.
- **Entraînement contradictoire** : technique d'apprentissage automatique utilisée pour rendre les modèles plus fiables. Tout d'abord, les développeurs construisent des « entrées contradictoires » (par exemple via le red-teaming) qui sont conçues pour faire échouer un modèle, et ensuite, ils entraînent le modèle à reconnaître et à gérer ce type d'entrées.
- **Apprentissage par renforcement à partir du retour d'information humain (RLHF)** : une technique d'apprentissage automatique dans laquelle un modèle d'IA est affiné en utilisant des évaluations ou des préférences fournies par l'homme comme signal de récompense, permettant au système d'apprendre et d'ajuster son comportement pour mieux s'aligner sur les valeurs et les intentions humaines grâce à un entraînement itératif.
- **Jailbreaking** : Génération et envoi d'invites conçues pour contourner les garde-fous et amener un système d'IA à produire du contenu nuisible, tel que des instructions pour la construction d'armes.

Les risques liés aux systèmes d'IA à usage général peuvent être atténués en partie en limitant leurs comportements. Par exemple, les décideurs politiques peuvent souhaiter empêcher les systèmes d'IA à usage général de fournir des informations dangereuses aux utilisateurs (par exemple sur la production d'armes ; voir [2.1.4. Attaques biologiques et chimiques](#)), d'être utilisés à des fins malveillantes (par exemple pour des cyberattaques ; voir [2.1.3. Cyberinfractions](#)) ou de connaître des dysfonctionnements entraînant des dommages (voir [2.2. Risques liés aux dysfonctionnements](#)). Le comportement d'un système est sûr s'il évite de telles erreurs, et un système est robuste s'il continue à se comporter de manière sûre dans un large éventail de circonstances. Au-delà de cela, un système est robuste sur le plan adversaire s'il maintient un comportement sûr même en présence d'un adversaire (par exemple un utilisateur humain) essayant de l'amener à effectuer des tâches nuisibles ou illégales.

Il existe des propositions sur la manière de construire des systèmes d'IA polyvalents dont le comportement est garanti en toute sécurité (1176), mais cela n'est pas possible sans avancées technologiques significatives et peut nécessiter des changements importants dans l'architecture des systèmes d'IA polyvalents actuels. La réglementation des systèmes actuels devra s'efforcer de garantir que leur formation et leur développement minimisent les risques de dysfonctionnements et d'utilisation abusive.

Depuis la publication du rapport intermédiaire, les attaquants et les défenseurs sont devenus plus efficaces pour tirer parti d'une compréhension plus approfondie du fonctionnement interne des systèmes d'IA afin d'induire ou de prévenir des comportements nuisibles, et l'avantage reste du côté des attaquants. De nouvelles méthodes pour résister aux attaques adverses en exploitant les concepts représentés en interne dans les réseaux neuronaux ont été développées à la fois pour les modèles d'images (1194*) et les modèles de langage (1195). Cependant, ces approches ne sont pas complètement robustes, et une autre étude récente a montré que les modèles de langage représentent en interne le refus des demandes nuisibles de manière simple.

Les attaques sont également faciles à exploiter (907). Dans l'ensemble, l'avantage reste généralement aux attaquants, qui peuvent inciter un modèle à adopter un comportement nuisible avec un effort modéré. Cependant, ces développements suggèrent que des recherches plus poussées sur les attaques et les défenses permettront probablement de progresser en matière d'interprétabilité. Si cela est vrai, de nouvelles avancées pourraient favoriser les défenseurs dans le cas de modèles à pondérations fermées, puisque les attaquants n'auront pas accès aux éléments internes du réseau neuronal dans ces cas.

Il est également apparu que les méthodes existantes de formation de modèles à usage général peuvent les amener à produire des résultats plus trompeurs (c'est-à-dire faux mais convaincants). Une étude récente a montré que dans le cas de questions particulièrement difficiles, la formation de systèmes d'IA à usage général pour maximiser l'approbation humaine des réponses a conduit les systèmes à masquer leurs erreurs et à les rendre plus difficiles à repérer par les humains, au lieu de devenir plus précis (608). D'autres études dans des environnements simulés ont révélé qu'une IA apprend à utiliser des stratégies néfastes (par exemple, cacher des informations ou exploiter les préjugés de son superviseur) pour recevoir un retour positif (1196) ou à modifier son environnement de formation pour augmenter sa récompense (599*), si l'IA dispose de suffisamment d'informations sur la manière de le faire. L'utilisation de l'IA pour aider les superviseurs à éviter les erreurs reste un problème difficile, mais des progrès modestes ont également été réalisés dans ce domaine, deux études récentes montrant des cas où les modèles deviennent plus faciles à superviser lorsqu'ils sont optimisés pour débattre entre eux (1197, 1198). Ces évolutions soulignent la nécessité de poursuivre les recherches sur les comportements encouragés par les méthodes de formation actuelles et de développer de nouvelles méthodes de formation qui offrent de meilleures incitations et des résultats généralement plus fiables par conception.

Les principales lacunes en matière de données probantes concernant la formation de modèles fiables sont les suivantes :

- Malgré les progrès récents (1010*, 1012, 1199), il n'est toujours pas certain que les méthodes d'interprétabilité, qui aident les chercheurs et les évaluateurs à comprendre le fonctionnement interne des modèles, soient suffisamment utiles pour éclairer de manière substantielle la formation et les tests des modèles. Il existe des études préliminaires à ce sujet (1076, 1200, 1201).
- Il n'est pas certain que les protocoles de « surveillance évolutive », dans lesquels les systèmes d'IA peuvent aider les humains, soient efficaces. évaluer leurs résultats peut fournir un levier puissant par lequel les modèles peuvent être formés pour être plus fiables même sur des problèmes difficiles (609*).
- Il n'existe actuellement aucune approche technique viable pour quantifier rigoureusement le risque de défaillances imprévues ou inattendues dans les grands systèmes d'IA à usage général. Bien que des recherches soient en cours sur l'obtention de garanties de sécurité probabilistes, il n'existe pas encore de technique pratique pour obtenir des garanties même approximatives.

Pour les décideurs politiques, les principaux défis sont les suivants :

- La recherche évolue très rapidement dans le domaine de la formation de l'IA, ce qui en fait une cible mouvante pour la réglementation.
- Il est difficile de quantifier le risque de modes de défaillance inattendus et imprévus. En outre, il est on ne sait pas clairement quelles sont les meilleures pratiques que les développeurs d'IA devraient utiliser pour détecter, répondre et atténuer les défaillances nouvellement découvertes afin de minimiser les risques.

Robustesse

Encourager un comportement sûr et correct pendant la formation au système

Il est difficile de spécifier précisément les objectifs des systèmes d'IA à usage général de manière à ne pas encourager involontairement des comportements nuisibles. Actuellement, les chercheurs ne savent pas comment spécifier des préférences et des valeurs humaines abstraites (telles que rapporter la vérité, comprendre et faire ce que veut un utilisateur ou éviter des actions nuisibles) d'une manière qui puisse être utilisée pour former des systèmes d'IA à usage général. De plus, étant donné les relations sociotechniques complexes intégrées dans les systèmes d'IA à usage général, il n'est pas certain qu'une telle spécification soit même possible. Après une phase initiale de pré-formation, les systèmes d'IA à usage général ont appris à imiter le comportement humain et sont ensuite généralement réglés pour optimiser des objectifs qui sont des proxys imparfaits des véritables objectifs du développeur (1031). Par exemple, les chatbots d'IA sont souvent réglés pour produire du texte qui sera évalué positivement par des évaluateurs humains, mais l'approbation de l'utilisateur est un proxy imparfait des avantages pour l'utilisateur. Des recherches ont montré que plusieurs chatbots largement utilisés font parfois correspondre leurs opinions déclarées à celles d'un utilisateur, sans tenir compte de la vérité (98, 522), ce qui peut créer des « chambres d'écho », et que la formation de systèmes d'IA à usage général pour satisfaire aux évaluations des évaluateurs humains peut inciter le système à fournir des réponses plus difficiles à vérifier qui masquent les erreurs du système (608). Il s'agit d'un défi permanent pour les systèmes d'IA à usage général (607, 1029, 1031, 1202*).

Les chercheurs disposent de méthodes pour mesurer si la formation incite au bon comportement en utilisant des expériences avec des évaluateurs humains, mais les résultats actuels sont préliminaires. Les expériences de « supervision évolutive » testent si un évaluateur peut diriger avec succès un système d'IA pour effectuer correctement une tâche que l'évaluateur est incapable de démontrer ou d'évaluer lui-même – par exemple, pour répondre à des questions (telles que des questions de sciences dures) qui nécessitent une expertise spécialisée pour être vérifiées (609*, 1203*). Cela fournit une vérification empirique solide que le protocole de formation utilisé incite au bon comportement. Les protocoles en cours de développement pour une supervision évolutive font souvent appel au système d'IA lui-même pour aider l'évaluateur, par exemple en le faisant s'engager dans un débat avec lui-même sur la bonne réponse (611*), et en laissant un évaluateur humain diriger le modèle sur la base de ce débat. Des expériences récentes de débat entre humains et IA montrent que cela peut améliorer la capacité des évaluateurs humains à déterminer les bonnes réponses à des questions difficiles (615*, 1198, 1204*), et les résultats préliminaires montrent que cela peut se traduire par une meilleure incitation à la formation (1197). Cependant, des résultats positifs n'ont été observés que sur une simple tâche de compréhension de lecture, avec des résultats mitigés pour d'autres tâches telles que des problèmes de mathématiques (1198). Ces méthodes n'ont pas été utilisées pour former une IA à usage général

systemes, mais les progrès dans ce domaine se poursuivent et des expériences de surveillance évolutives pourraient à un moment donné constituer un moyen pratique de mesurer la fiabilité avec laquelle les techniques de formation incitent au bon comportement.

Certains chercheurs travaillent sur des approches « sûres par conception » qui pourraient être en mesure de fournir des garanties de sécurité quantitatives. Au-delà de garantir que le processus de formation d'une IA code l'incitation à être sûre, il pourrait être possible de concevoir des systèmes d'IA qui garantissent quantitativement certains niveaux de sécurité (1176). Ces propositions reposent souvent sur une combinaison de trois éléments : premièrement, une spécification des résultats souhaités et indésirables (qui dans certains cas pourrait être une description en langage naturel des comportements souhaités et inacceptables), deuxièmement, un « modèle mondial » qui comprend la capture des relations de cause à effet (approximatives) et prédit les résultats des actions possibles que le système d'IA pourrait entreprendre, et troisièmement, un vérificateur qui vérifie si une action candidate donnée conduirait à des résultats prévus indésirables. L'objectif de ce processus est de garantir que des actions dangereuses ne soient pas entreprises. Si le modèle mondial capture des connaissances scientifiques, il s'appuiera généralement sur des hybrides « neuro-symboliques » d'IA à usage général et de techniques classiques utilisant les mathématiques formelles.

Français L'avantage des garanties et limites mathématiques est qu'elles peuvent fournir des assurances de sécurité même en dehors du domaine dans lequel l'IA a été formée et testée, contrairement aux vérifications ponctuelles et aux améliorations par essais et erreurs qui sont actuellement la norme pour évaluer et former des modèles d'IA à usage général. Cette approche explicite basée sur un modèle offre deux avantages supplémentaires : tout d'abord, parce qu'elle utilise une logique formelle et des lois de probabilité pour analyser des composants de connaissances clairement définis, ses conclusions sont plus fiables, compréhensibles et vérifiables que celles des systèmes d'IA traditionnels. Deuxièmement, elle permet de construire des systèmes d'IA non agentiques (non autonomes) qui peuvent faire progresser la science et les connaissances humaines tout en restant faciles à contrôler, évitant les risques potentiels qui accompagnent l'IA hautement agentique avancée (voir [2.2.3. Perte de contrôle](#)). Actuellement, cependant, des garanties de sécurité utiles et démontrables dans la pratique n'ont pas encore été démontrées pour les modèles et méthodes d'IA à usage général, et de nombreuses questions restent ouvertes afin d'atteindre ces objectifs pour les systèmes d'IA à grande échelle (1205).

Maintenir la qualité de la supervision humaine et de l'évaluation du comportement de l'IA

Les techniques de formation et d'évaluation de pointe reposent sur des retours d'information ou des démonstrations de la part d'êtres humains et sont donc limitées par les erreurs et les biais humains. Les développeurs peaufinent les systèmes d'IA polyvalents de pointe en faisant appel à une grande implication humaine. En pratique, cela implique des techniques qui exploitent des exemples d'actions souhaitées générés par des êtres humains (28) ou des retours d'information générés par l'homme sur des exemples issus de modèles (29, 30, 31*, 1182). Cela se fait à grande échelle, ce qui rend l'opération coûteuse et exigeante en main-d'œuvre. Cependant, l'attention, la compréhension et la fiabilité humaines ne sont pas parfaites (1182), ce qui limite la qualité des systèmes d'IA polyvalents qui en résultent (1206, 1207*, 1208). Même de légères imperfections dans les retours d'information des humains peuvent être amplifiées lorsqu'elles sont utilisées pour entraîner des systèmes hautement performants, avec des conséquences potentiellement graves (voir par exemple [2.2.3. Perte de contrôle](#)).

L'amélioration de la qualité et de la quantité de la supervision humaine peut contribuer à former des modèles plus robustes. Certaines recherches ont montré que l'utilisation de formes de feedback plus riches et plus détaillées de la part des humains peut permettre une meilleure supervision des modèles d'IA, mais au prix d'une augmentation du temps et des efforts de collecte des données (1209*, 1210, 1211). Pour collecter des ensembles de données plus volumineux, l'exploitation de systèmes d'IA à usage général pour automatiser partiellement le processus de feedback peut augmenter considérablement le volume de données (33*, 256*). Cependant, dans la pratique, la quantité de supervision humaine explicite utilisée pendant le réglage fin est très faible par rapport aux milliards de points de données utilisés dans la pré-formation sur les données Internet, et la supervision humaine peut donc être incapable de supprimer complètement les connaissances ou les capacités nuisibles de la pré-formation. L'amélioration des données de feedback de réglage fin ne constituera probablement qu'une partie de la solution à la robustesse coop

Améliorer la factualité des résultats du modèle

L'hallucination de faussetés est un défi, mais elle peut être réduite. En IA, « hallucination » fait référence à la propension des systèmes d'IA à usage général à produire des faussetés et du contenu inventé. Par exemple, les modèles de langage hallucinent généralement des citations, des biographies et des faits inexistantes (101, 102*, 103, 104, 105), ce qui pourrait poser des problèmes juridiques et éthiques impliquant la diffusion de fausses informations (1212). Il est possible mais difficile de réduire la tendance des systèmes d'IA à usage général à halluciner des résultats faux. Affiner explicitement les modèles d'IA à usage général pour les rendre plus véridiques - à la fois dans l'exactitude de leurs réponses et dans l'analyse de leur propre compétence -

est une approche pour relever ce défi (1213*). De plus, permettre aux modèles de langage d'accéder aux bases de données de connaissances lorsqu'on leur demande d'effectuer des tâches contribue à améliorer la fiabilité de leurs générations (838, 1214). D'autres approches détectent les hallucinations et informent l'utilisateur si le résultat généré n'est pas fiable (1215), effectuent des vérifications précises sur les déclarations individuelles faites par un modèle (1216) ou quantifient la confiance du modèle (1217). Cependant, la réduction des hallucinations reste un domaine de recherche très actif.

Améliorer la robustesse face aux pannes inattendues

Il est très difficile de garantir que les systèmes d'IA à usage général apprennent des comportements bénéfiques qui se traduisent de leurs contextes de formation à des contextes de déploiement réels à enjeux élevés. Parfois, des entrées inconnues qu'un système d'IA à usage général rencontre lors du déploiement peuvent provoquer des échecs inattendus (1218).

Tout comme les systèmes d'IA à usage général sont formés pour optimiser des objectifs proxy imparfaits, le contexte de formation peut également ne pas représenter de manière adéquate les situations réelles que les systèmes rencontreront après leur déploiement. Dans de tels cas, les systèmes d'IA à usage général peuvent toujours prendre des mesures nuisibles même s'ils sont formés avec un retour d'information correct fourni par l'homme (616, 1032, 1033). Par exemple, certains chercheurs ont découvert que les chatbots sont plus susceptibles de prendre des mesures nuisibles dans des langues sous-représentées dans leurs données de formation (1034). Une façon d'atténuer ces échecs consiste à utiliser des cadres d'évaluation qui testent de nombreuses combinaisons de conditions de déploiement, comme le cadre d'évaluation holistique des modèles linguistiques (HELM (1150)), qui énumère et teste des combinaisons de nombreuses tâches, profils d'utilisateur et langues différents, entre autres fonctionnalités. Une autre solution consiste à développer des méthodes par lesquelles les modèles peuvent estimer et communiquer leur incertitude dans de rares cas afin d'anticiper les erreurs (1219*, 1220*). Cependant, dans

En général, il est probablement impossible d'énumérer toutes les situations réelles possibles à des fins d'évaluation ou d'anticiper toutes les erreurs potentielles.

La compréhension des calculs internes d'un modèle peut aider les chercheurs à déterminer s'ils ont appris des solutions robustes. Il existe des méthodes permettant d'identifier automatiquement les caractéristiques (c'est-à-dire les modèles mathématiques) à l'intérieur d'un modèle de réseau neuronal qui correspondent à des concepts interprétables par l'homme (1009, 1013*, 1221, 1222*), notamment des personnes et des lieux spécifiques ainsi que des concepts et des comportements abstraits tels que des erreurs de code, la non-conformité à certaines opinions politiques ou des descriptions de la manière de créer des médicaments (1012). Ces caractéristiques peuvent servir de guide pour identifier les comportements dangereux ou indésirables dans les données d'entraînement d'un système ou ses résultats à une échelle plus grande que ce qui serait possible avec un examen humain seul. Les chercheurs ont tenté d'automatiser cet examen à l'aide d'un « agent d'interprétabilité automatisé » qui a accès à des outils d'interprétabilité. Une étude préliminaire montre que cela est possible à petite échelle (1201), et qu'il n'existe aucun obstacle évident à l'intensification de ce type de travail.

Des progrès récents ont été réalisés pour utiliser la compréhension du fonctionnement interne d'un modèle afin d'améliorer son comportement, mais cette approche nécessite davantage de travail. Malgré la difficulté de comprendre le fonctionnement interne des modèles, certaines techniques peuvent être utilisées pour guider des modifications spécifiques. Par rapport au réglage fin, ces méthodes peuvent parfois être des moyens plus efficaces en termes de calcul ou de données pour modifier la fonctionnalité des modèles. Les chercheurs ont utilisé diverses méthodes pour cela, basées sur des modifications des paramètres internes des modèles appris pendant l'entraînement (1223, 1224, 1225, 1226, 1227), des neurones (1221, 1228, 1229) ou des représentations (1199, 1230, 1231, 1232, 1233). Ces techniques sont imparfaites (1023), généralement limitées à des types de comportements très spécifiques (1227) et introduisent généralement des effets secondaires imprévus sur le comportement du modèle (1234), mais elles restent un domaine de recherche actif. On ne sait pas dans quelle mesure les méthodes actuelles offrent un moyen « utile et fiable » de comprendre et de concevoir des modèles d'IA à usage général (1026*).

Robustesse contradictoire : éviter les abus de modèles

Les utilisateurs de systèmes d'IA à usage général peuvent souvent contourner leurs mesures de protection grâce à des « jailbreaks » qui les incitent à se conformer à des demandes nuisibles. Même si un système se comporte toujours bien dans des conditions normales d'utilisation, un individu motivé peut toujours construire des entrées inhabituelles spécifiquement conçues pour faire échouer un système ou adopter des comportements indésirables (par exemple nuisibles) (1054). Les modèles de langage en particulier sont sujets à des « jailbreaks » à usage général qui peuvent les rendre beaucoup plus susceptibles de se conformer à des demandes nuisibles. Voici quelques exemples de méthodes de jailbreaking : inciter un système d'IA à adopter la personnalité de quelqu'un qui dirait le contenu nuisible (1053), l'amorcer avec des exemples de réponses nuisibles (1235*) ou faire des demandes dans une langue qui était rare dans les données d'entraînement du système (1236), ce qui pourrait accroître la vulnérabilité des modèles dans certains pays à revenu faible ou intermédiaire (PRFI) (voir le tableau 3.3 pour quelques exemples de jailbreaks). Bien que les jailbreaks puissent être partiellement évités après leur découverte, il est difficile de les anticiper lors du développement du modèle et, à l'heure actuelle, il est généralement facile de trouver de nouveaux jailbreaks qui fonctionnent pour les modèles de pointe. Cela étant, on ne sait pas dans quelle mesure les jailbreaks sont utilisés pour provoquer des comportements nuisibles par les systèmes d'IA en dehors d'un cadre de recherche.

Entraîner des modèles pour détecter et refuser les demandes nuisibles des adversaires

L'entraînement antagoniste permet d'améliorer la robustesse des systèmes d'IA de pointe, mais seulement dans une mesure limitée. L'« entraînement antagoniste » consiste d'abord à construire des « attaques » conçues pour faire réagir un modèle de manière indésirable, puis à entraîner le système à gérer ces attaques de manière appropriée. Les attaques contre les systèmes d'IA peuvent prendre de nombreuses formes et peuvent être générées par des humains ou des algorithmes. Une fois qu'une attaque antagoniste a été produite, l'entraînement sur ces exemples peut se dérouler comme d'habitude. L'entraînement contradictoire est devenu une technique couramment utilisée pour rendre les modèles plus robustes aux pannes, et est utilisé dans le développement des principaux systèmes d'IA à usage général (4*, 48*, 147*, 1158*, 1163, 1241). Cependant, il n'est pas suffisant en soi, car les systèmes entraînés de manière contradictoire restent généralement vulnérables aux attaques, en particulier avec des entrées multimodales (par exemple avec des images). De plus, la pertinence ou la nocivité potentielle des sorties d'un système d'IA ne peut pas toujours être évaluée en dehors du contexte dans lequel il est utilisé, ce qui n'est pas disponible lors de l'entraînement contradictoire (1242).

Rendre les systèmes d'IA à usage général plus résistants aux attaques imprévues est un problème ouvert et difficile, mais il existe des méthodes potentiellement prometteuses pour minimiser les dommages concernés.

L'entraînement contradictoire nécessite généralement des exemples précis d'échecs (598*, 1243). Ces limitations ont donné lieu à des jeux permanents de « chat et de la souris » dans lesquels certains développeurs mettent continuellement à jour les modèles en réponse aux vulnérabilités nouvellement découvertes. Le processus de recherche de vulnérabilités et de tentative d'induire un comportement indésirable est connu sous le nom de « red-teaming ». Une solution partielle à la vulnérabilité continue des modèles consiste simplement à produire et à s'entraîner sur des exemples plus contradictoires. Les méthodes automatisées de génération d'attaques peuvent aider à accroître la formation des adversaires (522, 904*, 1157, 1244). Cependant, le nombre exponentiellement élevé d'entrées possibles pour les systèmes d'IA à usage général rend difficile la recherche approfondie de tous les types d'attaques. Les méthodes d'interprétabilité pourraient être utiles ici (907), et des progrès préliminaires ont été réalisés pour améliorer la robustesse grâce à des méthodes qui fonctionnent sur les états internes du modèle (1076, 1195, 1200). Même si toutes les attaques ne peuvent pas être empêchées à l'avance, si elles peuvent être détectées rapidement au moment de l'exécution, les systèmes peuvent être efficacement adaptés pour s'en défendre : dans une étude, un système a eu plus de 95 % de succès pour se défendre contre les attaques après avoir vu un seul exemple du même type d'attaque (1245). Bien que la recherche sur ces mesures d'atténuation soit préliminaire, il est essentiel d'exiger une surveillance en direct, une réponse et des mesures d'atténuation de la formation des adversaires sur des systèmes d'IA potentiellement dangereux pour réduire les dommages causés par une mauvaise utilisation.

Les méthodes de « désapprentissage automatique » visent à supprimer certaines capacités indésirables des systèmes d'IA à usage général, mais les techniques actuelles suppriment souvent plutôt que de supprimer complètement ces capacités. Par exemple, le désapprentissage automatique peut supprimer certaines capacités qui pourraient aider les utilisateurs malveillants à fabriquer des explosifs, des armes biologiques, des armes chimiques et des cyberattaques (392). Le désapprentissage comme moyen de nier l'influence des données d'entraînement indésirables a été proposé à l'origine comme un moyen de protéger la vie privée et le droit d'auteur (821), abordé dans [2.3.6. Risques de violation du droit d'auteur](#). Les méthodes de désapprentissage pour supprimer les capacités dangereuses (892, 1246) incluent des méthodes basées sur le [réglage fin](#) (893*) et la [modification du fonctionnement interne](#) des modèles (392). Idéalement, le désapprentissage devrait rendre un modèle incapable de présenter le comportement indésirable même lorsqu'il est soumis à des attaques d'extraction de connaissances, à des situations nouvelles (par exemple des demandes dans différentes langues) ou à de petites quantités de réglage fin. Cependant, les méthodes actuelles de désapprentissage devraient permettre à un modèle d'être incapable de présenter le comportement indésirable, même lorsqu'il est soumis à des attaques d'extraction de connaissances, à des situations nouvelles (par exemple des demandes dans différentes langues) ou à de petites quantités

Les méthodes de désapprentissage suppriment souvent les informations nuisibles sans les supprimer de manière efficace (1247). Cela crée des défis pour la gouvernance, car les modèles peuvent sembler dépourvus de capacités nuisibles alors que celles-ci sont en fait simplement cachées et peuvent être réactivées. Les méthodes de désapprentissage actuelles peuvent également introduire des effets secondaires indésirables sur les connaissances souhaitables du modèle (1247). Il n'est pas certain que le désapprentissage d'une compétence nuisible puisse supprimer complètement la capacité du modèle à effectuer une tâche nuisible en combinant les compétences et les connaissances souhaitables. Le désapprentissage reste un domaine de recherche actif.

3.4.2. Suivi et intervention

INFORMATIONS CLÉS

- La surveillance et l'intervention sont des approches complémentaires pour prévenir l'IA dysfonctionnements et utilisations malveillantes. Les moniteurs inspectent les entrées et sorties du système, l'état du matériel, les composants internes du modèle et les impacts réels pendant l'utilisation des systèmes, déclenchant des interventions qui bloquent les actions potentiellement nuisibles. Les outils actuels peuvent détecter le contenu généré par l'IA, suivre le comportement du système et identifier les tendances préoccupantes dans ces cibles de surveillance. Cependant, les utilisateurs moyennement qualifiés peuvent souvent contourner ces mesures de protection par divers moyens techniques.
- Les méthodes d'interprétation et d'explication des modèles peuvent aider à surveiller les décisions de l'IA, mais les méthodes actuelles peuvent également produire des informations trompeuses. Les approches techniques permettant d'expliquer les résultats des systèmes d'IA aident les développeurs et les déployeurs à examiner minutieusement la prise de décision, bien que des études indiquent que ces méthodes peuvent produire des explications inexactes ou trop simplifiées du comportement complexe des modèles.
- Plusieurs niveaux de surveillance et d'intervention renforcent la protection contre les dysfonctionnements et les utilisations malveillantes. La combinaison de capacités de surveillance et d'intervention techniques avec l'intervention humaine permet de renforcer les garanties, même si ces mesures peuvent entraîner des coûts et des retards.
- Ces derniers mois, des progrès ont été réalisés en matière d'interprétabilité des modèles et de mesures de surveillance basées sur le matériel. Depuis la publication du rapport intermédiaire (mai 2024), la recherche sur l'interprétabilité des modèles a progressé pour commencer à expliquer les comportements des modèles, et les premiers travaux sur la surveillance basée sur le matériel préservant la confidentialité ont le potentiel d'améliorer la visibilité réglementaire du développement de l'IA.
- Les principaux défis pour les décideurs politiques sont de trouver un équilibre entre les mesures de sécurité et leurs coûts pratiques. Si les mesures de sécurité à plusieurs niveaux offrent une meilleure protection, elles entraînent également des retards opérationnels, soulèvent des problèmes de confidentialité et augmentent les coûts de déploiement. Les décideurs politiques doivent donc évaluer les exigences de sécurité par rapport à ces contraintes pratiques, en particulier compte tenu du décalage potentiel entre les mesures de sécurité et les incitations commerciales.

Définitions clés

- **Modèle** : programme informatique, souvent basé sur l'apprentissage automatique, conçu pour traiter des entrées et générer des sorties. Les modèles d'IA peuvent effectuer des tâches telles que la prédiction, la classification, la prise de décision ou la génération, constituant ainsi le cœur des applications d'IA.
- **Système** : une configuration intégrée qui combine un ou plusieurs modèles d'IA avec d'autres composants, tels que des interfaces utilisateur ou des filtres de contenu, pour produire une application avec laquelle les utilisateurs peuvent interagir.

- **Interprétabilité** : degré auquel les humains peuvent comprendre le fonctionnement interne d'un modèle d'IA, notamment pourquoi il a généré un résultat ou une décision particulière. Un modèle est hautement interprétable si ses processus mathématiques peuvent être traduits en concepts qui permettent aux humains de retracer les facteurs et la logique spécifiques qui ont influencé le résultat du modèle.
- **Contenu factice généré par l'IA** : contenu audio, texte ou visuel, produit par l'IA générative, qui représente des personnes ou des événements d'une manière qui diffère de la réalité, de manière malveillante ou trompeuse, par exemple en montrant des personnes faisant des choses qu'elles n'ont pas faites, en disant des choses qu'elles n'ont pas dites, en changeant le lieu d'événements réels ou en représentant des événements qui ne se sont pas produits.
- **Deepfake** : un type de faux contenu généré par l'IA, composé de contenu audio ou visuel, qui déforme les faits et présente des personnes réelles comme faisant ou disant quelque chose qu'elles n'ont pas réellement fait ou dit.
- **Informatique légale** : processus de traçage de l'origine et de la diffusion des médias numériques.
- **Filigane** : un motif subtil, souvent imperceptible, intégré dans un contenu généré par l'IA (comme du texte, des images ou de l'audio) pour indiquer son origine artificielle, vérifier sa source ou détecter une éventuelle utilisation abusive.
- **Défense en profondeur** : une stratégie qui comprend la superposition de plusieurs mesures d'atténuation des risques cas où aucune méthode unique existante ne peut assurer la sécurité.
- **L'humain dans la boucle** : une exigence selon laquelle les humains doivent superviser et approuver les processus automatisés dans les domaines critiques.
- **Agent IA** : une IA à usage général qui peut élaborer des plans pour atteindre des objectifs, effectuer de manière adaptative des tâches impliquant plusieurs étapes et des résultats incertains en cours de route, et interagir avec son environnement (par exemple en créant des fichiers, en effectuant des actions sur le Web ou en déléguant des tâches à d'autres agents) avec peu ou pas de surveillance humaine.

Les stratégies de surveillance et d'intervention sont appliquées aux systèmes d'IA (le package de déploiement complet qui comprend à la fois le modèle d'IA et des composants de sécurité supplémentaires), laissant le modèle inchangé. Contrairement aux stratégies décrites dans [la section 3.4.1. Entraînant des modèles plus fiables](#), les méthodes de surveillance et d'intervention sont intégrées au niveau du système et mises en œuvre dans le cadre du déploiement du système. Cette section décrit les stratégies de surveillance et d'intervention que les chercheurs et les développeurs utilisent pour les systèmes d'IA à usage général (voir la figure 3.2).

Les principales lacunes en matière de données probantes concernant la surveillance et l'intervention concernent la compréhension de l'efficacité des méthodes et de la facilité avec laquelle elles peuvent être contournées. Les techniques de surveillance et d'intervention sont, dans de nombreux cas, des mesures de protection simples et efficaces au niveau du système dans les cas d'utilisation typiques. Elles offrent une ligne de défense supplémentaire essentielle en plus des techniques au niveau du modèle évoquées dans [la section 3.4.1. Former des modèles plus fiables](#). De ce point de vue, il existe peu d'obstacles techniques à l'adoption généralisée de nombreuses techniques. Cependant, les scientifiques n'ont pas encore une compréhension quantitative approfondie de leur efficacité dans des contextes réels et de la facilité avec laquelle les méthodes de surveillance peuvent être coordonnées dans la chaîne d'approvisionnement de l'IA. L'un des principaux obstacles à l'adoption de techniques de surveillance et d'intervention très efficaces est de comprendre à quel point elles sont vulnérables au contournement actif par des utilisateurs malveillants.

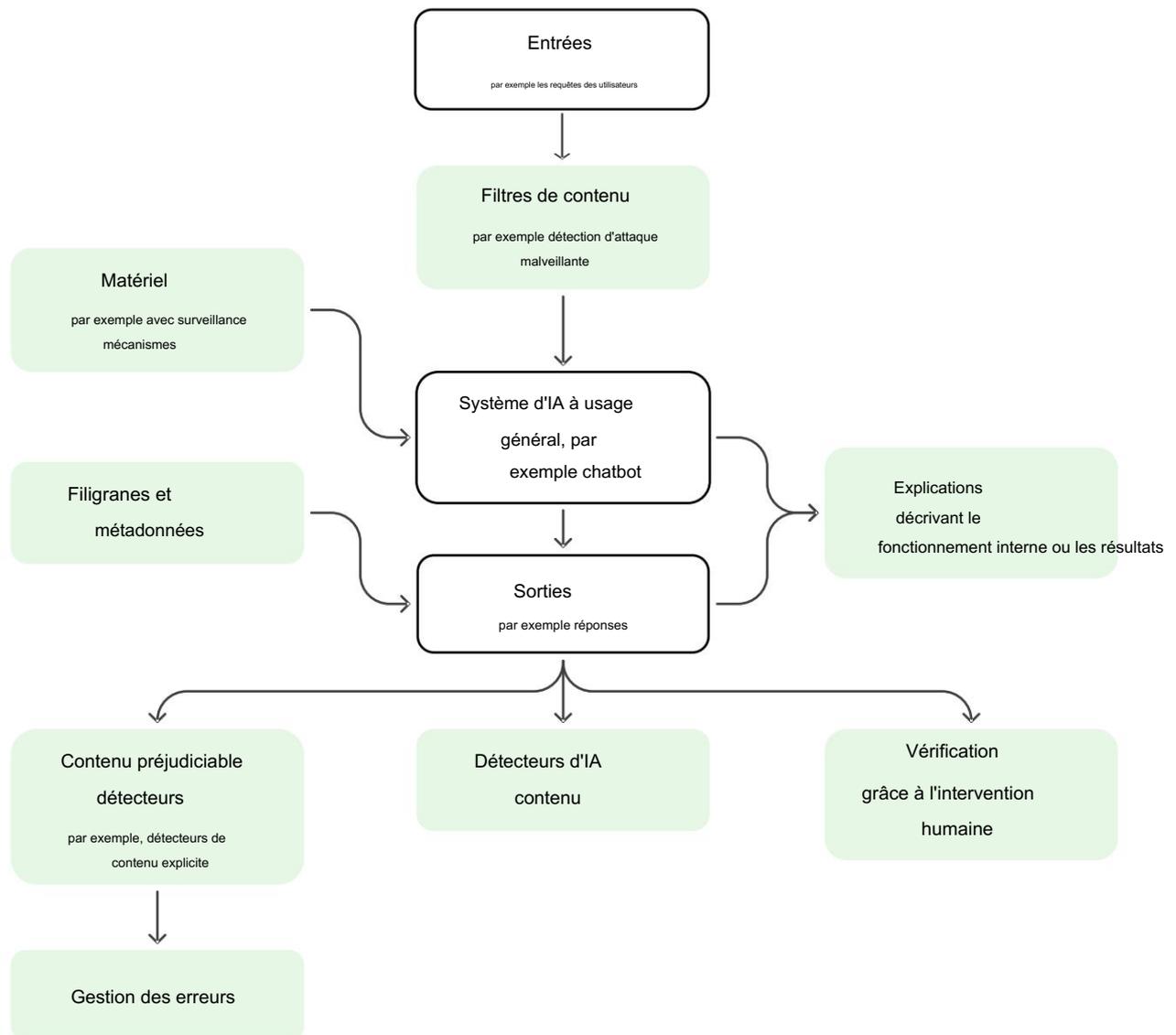


Figure 3.2 : Les techniques de surveillance et d'intervention sont des mesures de protection au niveau du système qui peuvent être appliquées aux entrées, aux sorties et aux modèles des systèmes d'IA à usage général afin d'aider les chercheurs et les développeurs à surveiller le comportement de l'IA et, si nécessaire, à intervenir. Source : International AI Safety Report.

Détection de contenu généré par l'IA

Le contenu généré par les systèmes d'IA à usage général – en particulier les « deepfakes » – pourrait avoir des effets néfastes à grande échelle (1248, 1249, 1250) (voir 2.1.1. [Dommages causés aux individus par le biais de faux contenus](#)). Toutefois, la capacité à distinguer les contenus authentiques des contenus générés par l'IA peut contribuer à réduire l'utilisation néfaste des modèles génératifs. Par exemple, si les navigateurs Web pouvaient afficher des avis de fiabilité sur les contenus susceptibles d'être générés par l'IA, cela contribuerait à lutter contre la propagation de fausses informations en ligne. Il existe une variété d'outils techniques permettant de détecter les contenus générés par l'IA. Aucun n'est parfait, mais ensemble, ils peuvent s'avérer extrêmement utiles pour la criminalistique numérique.

Il existe des techniques peu fiables mais néanmoins utiles pour détecter le contenu généré par l'IA. Tout comme les différents humains ont des styles artistiques et d'écriture discernables, les modèles d'IA génératifs le sont aussi. Certaines procédures ont été développées pour distinguer le texte généré par l'IA (332, 333, 337, 338, 1251, 1252, 1253) et les images (1254, 1255) du contenu généré par l'homme. Les méthodes de détection sont généralement basées soit sur des classificateurs spécialisés, soit sur l'évaluation de la probabilité qu'un exemple donné ait été généré par un modèle d'IA à usage général. Cependant, les méthodes existantes sont limitées et sujettes aux erreurs. Un défi important est que les systèmes d'IA à usage général ont tendance à mémoriser des exemples qui apparaissent dans leurs données d'entraînement. Pour cette raison, des extraits de texte courants (par exemple des documents historiques célèbres) ou des images d'objets courants (par exemple des œuvres d'art célèbres) sont parfois classés à tort comme étant générés par l'IA.

À mesure que le contenu généré par l'IA à usage général devient plus réaliste, il peut devenir de plus en plus difficile à détecter. Dans le même temps, les détecteurs de texte IA ont tendance à avoir des performances inégales selon les langues du monde, ce qui pose des problèmes d'égalité linguistique (1256).

Les « filigranes » – des motifs subtils mais distincts insérés dans les données générées par l'IA – facilitent la distinction du contenu généré par l'IA, mais ils peuvent être supprimés. Les filigranes sont des éléments souvent conçus pour être difficiles à remarquer pour un humain, mais faciles à identifier pour les algorithmes de détection.

Les filigranes se présentent généralement sous la forme de motifs imperceptibles insérés dans des pixels d'image ou de vidéo (290, 291, 292, 293, 294*, 1257), de signaux imperceptibles dans l'audio (295, 296) ou de biais stylistiques ou de choix de mots dans le texte (297, 1258, 1259, 1260, 1261). Les filigranes peuvent être utilisés pour détecter le contenu généré par l'IA avec une précision presque parfaite lorsqu'ils ne sont pas falsifiés. Comme indiqué dans [la section 2.1.1. Dompage causé aux individus par du faux contenu.](#) ils peuvent être utilisés pour détecter le faux contenu généré par l'IA. Ils constituent une stratégie imparfaite pour détecter le contenu généré par l'IA (en particulier le texte) car ils peuvent être supprimés par de simples modifications des données (298*, 299, 333, 1262). Cependant, cela ne signifie pas qu'ils ne sont pas utiles. Par analogie, les empreintes digitales sont faciles à éviter ou à supprimer, mais elles restent très utiles en criminalistique. Enfin, il existe des inquiétudes concernant la confidentialité et l'utilisation abusive potentielle de la technologie du tatouage numérique, car elle pourrait être utilisée pour suivre et identifier les utilisateurs (300).

Les filigranes peuvent également être utilisés pour indiquer un contenu authentique, non généré par l'IA. La certification de l'authenticité des données fait partie de la « provenance des données ». Contrairement à l'insertion de filigranes dans un contenu général généré par l'IA, une autre approche consiste à insérer automatiquement des filigranes dans un contenu non généré par l'IA (1263). Cependant, cela nécessitera souvent des modifications du matériel et du logiciel des appareils d'enregistrement physiques. Ces méthodes de provenance seraient très difficiles à falsifier au niveau de l'appareil. Certains chercheurs travaillent à l'élaboration de méthodes et de normes communes pour retracer l'origine des médias, y compris l'utilisation de méthodes de cryptage pour prouver l'authenticité qui sont difficiles à falsifier (par exemple CPPA (1264) ; AIMASC (1265)).

Les « métadonnées » et les journaux d'activité du système aident à la criminalistique numérique. La « criminalistique numérique » fait référence à la science de l'identification et de l'analyse des preuves numériques (1266, 1267, 1268, 1269, 1270). Il est courant que les données soient enregistrées avec des « métadonnées » qui donnent un contexte supplémentaire sur les données stockées. Ces métadonnées sont utiles (et couramment utilisées) pour retracer l'origine des données. Par exemple, de nombreux appareils mobiles enregistrent des fichiers image et audio à l'aide de la norme Exchangeable Image File Format (ExIF) (1271) qui peut stocker des informations sur les paramètres de l'appareil photo, l'heure, l'emplacement et d'autres détails.

Les métadonnées pourraient être utilisées pour aider à suivre les informations sur le fait que les données ont été générées par un système d'IA à usage général et, si tel est le cas, d'autres détails sur la manière dont cela a été fait. Par exemple, les développeurs et les déployeurs pourraient attacher des identifiants aux actions entreprises par un système d'IA (1272, 1273). Les développeurs et les déployeurs peuvent également enregistrer des « journaux d'activité » pour suivre le comportement du système, afin d'améliorer la surveillance au fil du temps (1272). En outre, le simple ajout d'étiquettes d'avertissement au contenu généré par l'IA peut contribuer à réduire la propagation de la désinformation. Une étude a révélé que ces étiquettes amélioraient la détection des deepfakes par les humains de 10,7 % à 21,6 % (289). Les métadonnées peuvent généralement être falsifiées, mais les preuves suggèrent que l'utilisation de signatures numériques cryptées peut permettre une preuve d'authenticité d'une manière très difficile à contrefaire (1274).

Au-delà des interventions techniques, des initiatives d'éducation aux médias numériques ont également été proposées pour lutter contre les faux contenus générés par l'IA (1275). Certaines études ont montré que les interventions d'éducation aux médias peuvent améliorer la capacité des participants à détecter les faux contenus (1276, 1277, 1278, 1279).

Cependant, en général, les données sur les effets des interventions en matière d'éducation aux médias numériques sont mitigées, en partie en raison des grandes variations dans les contextes d'étude et les modèles d'intervention (1279). Voir [2.1.1. Dommages causés aux individus par le biais de faux contenus](#) pour une discussion plus approfondie sur les faux contenus.

Détecter et se défendre contre les contenus nuisibles

Bien qu'il n'existe pas de mesure de sécurité parfaite, le fait de disposer de plusieurs niveaux de protection et de mesures de protection redondantes augmente la confiance dans la sécurité (une stratégie connue sous le nom de « défense en profondeur »). Bien que la présente section se concentre sur les approches techniques, les systèmes ne sont pas déployés dans le vide.

Les intégrer dans un système sociotechnique qui cherche à maintenir la sécurité et la performance est essentiel au processus continu d'identification, d'étude et de défense contre les dommages (également abordé au [point 3.1](#)).

[\(Aperçu de la gestion des risques\)](#) Cette section décrit diverses méthodes techniques complémentaires de détection et de défense contre les comportements nuisibles des systèmes d'IA à usage général.

La détection d'anomalies et de comportements potentiellement dangereux permet de prendre des précautions. Certaines méthodes ont été développées pour aider à détecter des entrées ou des comportements anormaux des systèmes d'IA (1280, 1281, 1282). Par exemple, les utilisateurs trompent parfois les modèles de langage pour qu'ils se comportent de manière nuisible en leur faisant encoder leurs réponses dans un texte chiffré (460, 1063*) qui ne ressemble pas du tout à du texte normal. Il est également parfois possible de détecter une proportion significative d'entrées (1243, 1283), d'états internes (1284, 1285, 1286*, 1287) ou de sorties (1287, 1288, 1289, 1290*, 1291)

impliqués dans des comportements nuisibles tels que l'assistance à des tâches dangereuses. Une fois détectés, les exemples risqués peuvent être envoyés à un processus de traitement des erreurs ou signalés pour une enquête plus approfondie. Par exemple, les données signalées comme nuisibles peuvent être bloquées par un filtre ou modifiées pour supprimer le contenu nuisible.

La présence d'un humain dans la boucle permet une surveillance directe et des commandes manuelles, mais peut être extrêmement coûteuse. Les humains dans la boucle sont coûteux par rapport aux systèmes automatisés. Cependant, lorsqu'il existe un risque élevé qu'un système d'IA à usage général prenne des mesures inacceptables, la présence d'un humain dans la boucle peut être essentielle. De même, les commandes manuelles sont standard dans les voitures dotées de modes de conduite autonome (1292). En attendant, les humains et les systèmes d'IA à usage général peuvent

Les décisions sont parfois prises de manière collaborative. Au lieu d'apprendre aux systèmes d'IA à usage général à agir au nom d'un humain, le paradigme de coopération homme-IA vise à combiner les compétences et les forces des systèmes d'IA à usage général et des humains (1293, 1294*, 1295, 1296, 1297, 1298, 1299).

Cependant, dans de nombreuses situations, il n'est pas pratique d'avoir un humain dans la boucle, par exemple lorsque la prise de décision est trop rapide (comme dans le cas d'applications de chat avec des millions d'utilisateurs), lorsque l'humain n'a pas suffisamment de connaissances du domaine ou lorsque les biais et les erreurs humaines peuvent exacerber les risques (1300). Les humains impliqués dans la boucle de prise de décision automatisée ont également tendance à présenter un « biais d'automatisation », ce qui signifie qu'ils accordent une plus grande confiance au système d'IA que prévu (1301). Dans les cas où l'intervention humaine n'est pas pratique, des approches hybrides impliquant un mélange de surveillance et d'intervention humaines et automatisées sont possibles.

Des protocoles de fonctionnement sécurisés peuvent être conçus pour des systèmes d'IA à usage général dotés de capacités potentiellement dangereuses. Les agents d'IA à usage général qui peuvent agir de manière autonome et sans limitation sur le Web ou dans le monde physique présentent des risques élevés (voir [3.2.1. Défis techniques pour la gestion des risques et l'élaboration des politiques](#) et [2.2.3. Perte de contrôle](#)). Pour les systèmes d'IA à usage général dotés de capacités potentiellement risquées, limiter les façons dont ils peuvent influencer directement le monde facilite leur surveillance et leur gestion (1302, 1303). Par exemple, si un système d'IA à usage général agentique a une capacité illimitée d'accéder aux systèmes de fichiers d'un ordinateur et/ou d'exécuter du code personnalisé, il est plus sûr d'exécuter cet agent dans un environnement informatique ad hoc que directement sur l'ordinateur de l'utilisateur (22*). Cependant, ces approches peuvent être difficiles à mettre en œuvre pour les applications dans lesquelles un système doit agir directement dans le monde. Dans ces cas, il est parfois difficile, même pour les humains, d'anticiper le moment où une action pourrait être nuisible.

Expliquer les actions du système d'IA

Certaines techniques peuvent être utilisées pour aider à expliquer pourquoi les systèmes d'IA à usage général déployés agissent comme ils le font. Comprendre pourquoi les systèmes d'IA à usage général agissent comme ils le font est utile pour évaluer les capacités, diagnostiquer les préjudices et déterminer la responsabilité en cas de préjudice (1304, 1305, 1306). Bien que cela puisse être utile, le simple fait de demander aux modèles de langage d'IA à usage général des explications sur leurs décisions peut également conduire à des réponses trompeuses (97, 1307). Pour augmenter la fiabilité des explications des modèles, les chercheurs travaillent sur des stratégies améliorées d'incitation et de formation (1308*, 1309*, 1310, 1311). Parallèlement, d'autres techniques d'explication des actions des modèles d'IA à usage général (1312, 1313) peuvent parfois aider à trouver des problèmes dans les modèles (1163). Cependant, expliquer correctement les actions des modèles d'IA à usage général est un problème difficile en raison de leur taille et de leur complexité. Certaines recherches visent à développer des techniques permettant d'aider les humains à interpréter les calculs des systèmes d'IA à usage général (1010*, 1011*, 1012). Les techniques permettant d'expliquer les décisions prises par les modèles sont reconnues comme un élément utile de la boîte à outils d'évaluation des modèles (1314).

Ces méthodes n'apportent toutefois qu'une compréhension partielle. Elles reposent sur des hypothèses importantes et des recherches supplémentaires sont nécessaires pour démontrer leur utilité dans la pratique.

Suivi et interventions avec du matériel spécialisé

Les mécanismes de surveillance préservant la confidentialité intégrés au matériel informatique apparaissent comme une alternative plus fiable et digne de confiance à la surveillance basée sur des logiciels ou à l'auto-déclaration. Le calcul est au cœur du développement et du déploiement de systèmes d'IA modernes à usage général, et la quantité de calcul utilisée pour la formation et l'inférence est corrélée à la capacité d'un système d'IA (voir [1.3. Capacités dans les années à venir](#)). La recherche sur les mécanismes matériels préservant la confidentialité vise à permettre aux décideurs politiques de surveiller et de vérifier certains aspects des systèmes d'IA à usage général pendant la formation et le déploiement, tels que l'utilisation du calcul, sans s'appuyer sur les rapports des développeurs d'IA. Par exemple, les recherches sur ces mécanismes soutiennent qu'ils permettent techniquement de vérifier les détails d'utilisation tels que l'heure et le lieu d'utilisation (1315, 1316), les types de modèles et de processus exécutés (1317, 1318), ou de fournir des preuves qu'un modèle particulier a été formé (1319, 1320). Si cela est possible, ces mécanismes peuvent être appliqués à de nombreuses questions de gouvernance, telles que la vérification du respect des accords internationaux, même au-delà des frontières (270). Certains pays envisagent de conclure des accords internationaux en raison des pressions concurrentielles entre les pays et de leur effet sur les incitations à gérer minutieusement les risques (voir [3.2.2. Défis sociétaux pour la gestion des risques et l'élaboration des politiques pour une analyse de cette dynamique](#)). Dans ce contexte, les pays peuvent résister au suivi et à la vérification des accords en raison de préoccupations concernant la propriété intellectuelle et les avantages concurrentiels. Des mécanismes de vérification basés sur le matériel sont parfois envisagés pour remédier à cette lacune, car ils pourraient permettre de surveiller des indicateurs clés tout en préservant la confidentialité des systèmes d'IA propriétaires et des données de formation. Cependant, ces applications en sont encore au stade des premières recherches (270).

Bien qu'une grande partie des fonctionnalités requises pour les mécanismes basés sur le matériel existent sur les puces d'IA actuelles, la surveillance basée sur le matériel n'a pas encore été prouvée à grande échelle et pourrait menacer les intérêts des utilisateurs si elle était mise en œuvre de manière aléatoire. Certains mécanismes basés sur le matériel sont largement déployés dans des contextes extérieurs à l'IA, comme les enclaves sécurisées d'Apple, qui permettent au fabricant de restreindre les applications installées sur ses appareils (1321*). Certaines puces d'IA de premier plan, comme l'unité de traitement graphique (GPU) H100, disposent déjà d'une partie du matériel nécessaire sous la forme de Confidential Computing (1322*). Néanmoins, certains mécanismes de surveillance et de vérification basés sur le matériel pour l'IA pourraient eux-mêmes être compromis par un attaquant bien doté en ressources, ce qui pourrait entraîner une fuite d'informations sensibles (1323).

3.4.3. Méthodes techniques de protection de la vie privée

INFORMATIONS CLÉS

- Les systèmes d'IA à usage général affectent la vie privée par la perte de confidentialité des données, le manque de transparence, traitement non autorisé des données et nouvelles formes d'abus. Ces risques sont décrits au [point 2.3.5. Risques pour la vie privée.](#)
- Plusieurs méthodes existent tout au long du cycle de vie de l'IA pour protéger la confidentialité. Il s'agit notamment de :
 - Depuis la publication du rapport intermédiaire (mai 2024), les méthodes de protection de la vie privée se sont développées pour répondre à l'utilisation croissante de l'IA dans des domaines sensibles. Cela inclut les assistants pour smartphone, les agents d'IA, les assistants vocaux toujours à l'écoute ou l'utilisation dans les soins de santé ou la pratique juridique. Il existe un intérêt croissant pour garantir la confidentialité et le consentement dans toutes ces utilisations, avec de nouvelles recherches et des mises en œuvre pratiques à l'appui. La suppression des informations personnelles identifiables (PII) et du contenu indésirable des données d'entraînement de l'IA à usage général, bien qu'encore difficile et incomplète, est un processus rentable, faisable et efficace pour réduire les risques. Des mécanismes conviviaux de contrôle et de traçage des données personnelles pourraient y contribuer.
- Les méthodes de protection de la vie privée dans le domaine de l'IA évoluent rapidement, ce qui crée des défis en matière de politiques publiques. Les méthodes visant à réduire les risques liés à l'IA à usage général sont complexes et continuent de se développer à un rythme soutenu, affectant de nombreux domaines de la chaîne d'approvisionnement et créant un environnement difficile pour l'élaboration des politiques publiques.

Définitions clés

- Confidentialité : droit d'une personne ou d'un groupe à contrôler la manière dont les autres accèdent à ses données sensibles ou les traitent. informations et activités.
- Informations personnelles identifiables (IPI) : toute donnée permettant d'identifier directement ou indirectement une personne (par exemple, un nom ou un numéro d'identification). Il s'agit d'informations pouvant être utilisées seules ou combinées à d'autres données pour identifier une personne de manière unique.
- Données sensibles : informations qui, si elles étaient divulguées ou mal traitées, pourraient entraîner des dommages, embarras, inconvénient ou injustice envers un individu ou une organisation.
- Minimisation des données : pratique consistant à collecter et à conserver uniquement les données directement nécessaires à un objectif spécifique, et à les supprimer une fois cet objectif atteint.
- Agent IA : une IA à usage général qui peut élaborer des plans pour atteindre des objectifs, effectuer de manière adaptative des tâches impliquant plusieurs étapes et des résultats incertains en cours de route, et interagir avec ses

environnement – par exemple en créant des fichiers, en effectuant des actions sur le Web ou en déléguant des tâches à d'autres agents – avec peu ou pas de surveillance humaine.

- Deepfake : un type de faux contenu généré par l'IA, composé de contenu audio ou visuel, qui déforme les faits et présente des personnes réelles comme faisant ou disant quelque chose qu'elles n'ont pas réellement fait ou dit.

Les méthodes et techniques d'atténuation des risques pour la vie privée liés à l'IA à usage général couvrent différentes catégories de risques.

[2.3.5. Les risques pour la vie privée](#) sont classés en trois grandes catégories : risques de formation (risques liés à la formation sur les données, en particulier les données sensibles) ; risques d'utilisation (risques liés à l'IA à usage général qui traite des informations sensibles pendant l'utilisation) ; et risques de préjudice intentionnel (risques liés aux acteurs malveillants appliquant l'IA à usage général pour porter atteinte à la vie privée des individus). Cette section examine les techniques d'atténuation pour chacune de ces catégories, en décrivant les techniques émergentes d'amélioration de la vie privée (1324) pour les catégories concernées. D'autres atteintes à la vie privée peuvent survenir à partir d'acteurs malveillants utilisant l'IA à usage général pour le harcèlement, les deepfakes non consentuels ou le vol d'informations sensibles ([2.1 Risques liés à une utilisation malveillante](#)), qui sont difficiles mais possibles à atténuer comme indiqué dans [3.4 Atténuation et surveillance des risques](#) et [2.1.3 Cyberinfraction](#)

Il est important et faisable de minimiser les informations personnellement identifiables dans les données de formation, mais cela représente un défi (Réduction des risques liés à la formation). L'IA à usage général est formée sur de grands ensembles de données collectés à partir de nombreuses sources, y compris le Web public. Ces données peuvent inclure des informations personnelles identifiables (1325, 1326), qui peuvent être reproduites lors de l'utilisation de modèles d'IA (827, 828, 1327, 1328). Les entreprises peuvent également utiliser leurs données propriétaires pour former des modèles (1329*). Les ensembles de données ouverts utilisés pour former l'IA à usage général tentent souvent de supprimer les informations personnelles identifiables (878, 1325) (bien que tous ne le fassent pas (1330)), mais peuvent manquer certaines informations personnelles identifiables. Sans normes plus claires pour la composition et l'inclusion éventuelle d'informations personnelles identifiables dans les ensembles de données (883, 1331), le nettoyage complet des données de formation pour l'IA à usage général à grande échelle sera difficile, mais le nettoyage des données reste un processus rentable, faisable et efficace en attendant de réduire les risques pour la vie privée.

La mise en œuvre de mécanismes conviviaux permettant aux individus de contrôler et de tracer leurs données, tels que des tableaux de bord pour la gestion des autorisations et des systèmes sécurisés de provenance des données, pourrait améliorer la transparence et la responsabilité dans les systèmes d'IA à usage général (Réduction des risques liés à la formation). Cela pourrait permettre aux individus de suivre la manière dont leurs données sont utilisées et partagées, d'établir des processus transparents pour que les individus puissent accéder à leurs données, les consulter, les corriger et les supprimer, ainsi que de suivre comment et où d'autres personnes profitent de leurs données (1332, 1333). Cela est possible pour les données détenues par l'utilisateur et, dans une moindre mesure, pour les données contenues sur les fournisseurs de services numériques (tels que les plateformes de médias sociaux) qui peuvent fournir des options de refus pour l'utilisation des données ou la formation (bien que les utilisateurs ne soient souvent pas conscients de leurs contributions à la formation de l'IA ou des risques de violation de la vie privée) (847, 1334, 1335). Si les données sont déjà accessibles au public sur le Web public, il est et restera beaucoup plus compliqué de contrôler la manière dont ces données sont utilisées pour l'IA à usage général.

Les approches de préservation de la confidentialité pour la formation sur des données sensibles sont limitées à l'IA à usage général (Réduction des risques de formation). Diverses techniques de confidentialité peuvent être appliquées aux modèles d'IA pour protéger la confidentialité des individus tout en permettant de tirer des informations utiles des données (1336, 1337).

Cependant, ces techniques peuvent considérablement altérer la précision du modèle (souvent appelée

Les techniques de « compromis vie privée-utilité » présentent des défis lorsqu'elles sont appliquées à des modèles de grande taille et peuvent ne pas convenir à tous les cas d'utilisation, en particulier pour les modèles d'IA à usage général formés sur du texte (1328). Pour les domaines contenant des données très sensibles (par exemple médicales ou financières), il peut être possible d'obtenir de solides garanties de confidentialité en adaptant de puissants modèles d'IA à usage général qui sont d'abord pré-entraînés sur des données accessibles au public sur Internet (1338*, 1339, 1340), mais de telles techniques ont rarement été appliquées dans la pratique jusqu'à présent. Une autre solution consiste à utiliser des données synthétiques (données, telles que du texte ou des images, qui ont été générées artificiellement, souvent par d'autres systèmes d'IA) pour éviter d'utiliser des données sensibles dans les pipelines de formation (1341*, 1342). Cependant, les chercheurs ont démontré qu'il existe un compromis important entre vie privée et utilité et qu'une forte confidentialité différentielle est toujours nécessaire pour la confidentialité (1343, 1344, 1345, 1346). La confidentialité différentielle fonctionne en ajoutant un bruit soigneusement calibré au processus de formation, limitant ainsi la quantité d'informations que le modèle peut apprendre sur les données d'un individu tout en lui permettant d'apprendre des modèles utiles à partir de l'ensemble de données. Si les données synthétiques sont très utiles, elles peuvent contenir autant d'informations que les données d'origine et permettre en grande partie les mêmes attaques (1347, 1348, 1349).

L'IA à usage général de capacité moyenne est de plus en plus capable de fonctionner entièrement sur des appareils grand public tels que les smartphones, ce qui permet aux utilisateurs d'utiliser l'IA à usage général sans envoyer de données personnelles à des serveurs externes (Réduction des risques d'utilisation). Alors que les systèmes d'IA à usage général les plus performants continueront d'être limités aux centres de données en raison de leur taille (156*), des systèmes d'IA plus petits qui peuvent répondre à des questions sur les données personnelles et effectuer des opérations téléphoniques de base au nom d'un utilisateur sont de plus en plus déployés sur des appareils grand public tels que les smartphones et autres appareils périphériques (4*, 37*, 841*). L'exécution de l'IA à usage général sur l'appareil signifie que les demandes de l'utilisateur et toutes les données personnelles auxquelles l'IA accède pour répondre à l'utilisateur n'ont pas besoin d'être envoyées à un serveur cloud externe, ce qui réduit les risques de fuite de données. Cependant, pour les tâches complexes, l'utilisation de l'IA à usage général est encore souvent externalisée vers (exécutée sur) des serveurs cloud, ce qui nécessite que les données personnelles et les demandes soient envoyées vers le cloud (via Internet).

Le déploiement sécurisé de systèmes d'IA à usage général dans le cloud est important lors de la manipulation de données sensibles (Réduction des risques d'utilisation). De nombreux grands modèles d'IA à usage général ne peuvent être exécutés que dans des centres de données, ce qui signifie que l'utilisation de données sensibles avec ces modèles nécessite l'envoi de ces données vers des emplacements externes. La sécurisation de ces déploiements est une tâche essentielle pour l'IA à usage général (844) et peut aider à empêcher la fuite d'informations privées. Les récents déploiements à grande échelle ont permis de créer des solutions de sécurité de bout en bout pour résoudre ce problème, mais des recherches supplémentaires sont nécessaires pour sécuriser ces déploiements (844).

Des approches cryptographiques solides pour exécuter l'IA de manière confidentielle et sécurisée de bout en bout existent, mais ne sont pas encore applicables à l'IA à usage général (Réduction des risques d'utilisation). Des recherches ont montré que de petits modèles d'IA peuvent être exécutés en combinaison avec des outils cryptographiques tels que le chiffrement homomorphe (1350), les preuves à connaissance nulle (1351), le calcul multipartite (1352, 1353) et les protections matérielles (telles que le calcul confidentiel sur les GPU NVIDIA H100) (1354, 1355, 1356*) pour permettre à la fois la confidentialité des entrées et la vérifiabilité du calcul sécurisé. Cependant, ces techniques imposent des coûts importants (les différentes méthodes peuvent différer de plusieurs ordres de grandeur) et n'ont pas été adaptées aux modèles les plus grands et les plus performants actuellement formés. 'Confidential

Le « calcul » avec les GPU H100 se distingue comme la seule approche cryptographique actuelle utilisable avec de grands modèles, mais ce n'est pas une solution complète pour le chiffrement ou la confidentialité de bout en bout.

Les avancées futures dans des domaines connexes pourraient permettre à ces techniques de sécurité renforcées de devenir pratiques pour l'IA à usage général à l'avenir (1177, 1357).

Les pratiques telles que la minimisation des données, la limitation des finalités et d'autres mesures de protection des données continueront d'être importantes avec l'IA à usage général, et les réglementations existantes en matière de confidentialité continueront de jouer un rôle dans la détermination de l'utilisation appropriée des données personnelles (Réduction des risques liés à la formation et à l'utilisation).

De nombreuses juridictions où l'IA à usage général sera utilisée disposent de réglementations existantes qui limitent ou établissent des lignes directrices sur la manière dont les données personnelles peuvent être utilisées (822, 1358). Dans de nombreux cas, les principes sous-jacents à ces réglementations s'appliquent déjà à la manière dont l'IA à usage général interagit avec les données personnelles ou sensibles et les utilise.

Les acteurs malveillants peuvent être en mesure d'utiliser l'IA à usage général pour violer la vie privée d'autrui par le biais d'un harcèlement amélioré par l'IA à usage général (Réduction des risques de préjudice intentionnel). Le contenu ci-dessus a principalement abordé les risques pour la vie privée liés à l'utilisation de données sensibles ou privées lors de la formation ou de l'utilisation de systèmes d'IA à usage général. Il existe également un risque distinct pour la vie privée lié aux acteurs malveillants utilisant l'IA à usage général pour améliorer les pratiques existantes de violation de la vie privée. L'IA à usage général peut déduire les attributs personnels des individus à un coût inférieur, à une vitesse plus élevée et à une plus grande échelle que les humains (483*, 846, 1047). Cela pourrait permettre, par exemple, à un acteur malveillant de rechercher dans de grandes violations de données et des informations publiques pour déduire les attributs des individus, déduire des informations sur le contenu public (comme l'endroit où une image a été prise) et effectuer des actions automatisées pour aider à l'exploitation de la vie privée, comme le phishing personnalisé automatique ou le harcèlement ciblé activé par l'IA à usage général. Certains cadres juridiques visent à tenir les créateurs et les distributeurs responsables des utilisations malveillantes (1359) et à fournir des recours aux personnes dont la vie privée a été violée.

D'autres capacités d'IA à usage général, telles que les attaques de cybersécurité avancées visant à extraire des informations privées ou les deepfakes non consentis, peuvent également aggraver cette tendance. Ces résultats pourraient être en partie évités grâce à des mesures techniques d'atténuation améliorées et sont similaires aux problèmes décrits ailleurs dans les sections [3.4 Atténuation et surveillance des risques](#) et [2.1 Risques liés à une utilisation malveillante](#).

Les systèmes d'IA à usage général peuvent également améliorer la confidentialité en soutenant les pratiques de cybersécurité dans le développement et en expliquant les risques aux utilisateurs. Bien que l'IA à usage général crée de nombreux risques pour la confidentialité, elle peut également contribuer à les atténuer. L'IA à usage général peut être utilisée dans les plateformes et outils de développement de logiciels, qui peuvent aider les développeurs à concevoir des logiciels sécurisés et à analyser les bases de code pour détecter d'éventuelles failles de sécurité (1047) (voir [2.1.3. Cyberinfraction pour en savoir plus sur l'utilisation de systèmes d'IA à usage général pour corriger les vulnérabilités logicielles](#)). Pour les utilisateurs, comprendre les risques pour la confidentialité et surveiller l'exposition personnelle est un défi. La narration et les explications centrées sur l'utilisateur des risques et des stratégies de sécurité personnelle en ligne sont importantes (1360) et pourraient être communiquées à l'aide de systèmes d'IA à usage général. Les systèmes d'IA pourraient également être utilisés pour aider à suivre où les données personnelles sont utilisées et à communiquer ces résultats aux utilisateurs.

Depuis la publication du rapport intermédiaire (mai 2024), des efforts accrus ont été déployés pour améliorer la qualité des données utilisées pour former l'IA à usage général, renforcer la sécurité matérielle du déploiement des systèmes d'IA et permettre l'exécution et le stockage de modèles localement sur des appareils personnels. Alors que l'IA à usage général devient de plus en plus accessible sur des appareils personnels tels que les assistants sur les smartphones (841*) et dans des contextes sensibles tels que les soins de santé (1361*), les outils de sécurité robustes pour héberger l'IA à usage général avec des garanties de confidentialité vérifiables deviennent plus courants (1362). Cette sécurité améliorée lors du déploiement (à la fois sur l'appareil et dans le cloud) est complétée par des travaux de filtrage des PII à partir des données de pré-formation à l'échelle du Web (878).

La capacité récente des systèmes d'IA à usage général à agir et à planifier de manière autonome au nom des utilisateurs (en tant qu'agents d'IA) a entraîné de nouveaux risques pour la vie privée (673, 1363).

D'autres considérations relatives à la confidentialité en aval sont également importantes. Par exemple, un certain nombre d'experts ont averti que si les agents d'IA devenaient impossibles à distinguer des humains authentiques sur le Web, la lutte contre ces systèmes entraînerait des risques d'identification massive (et de surveillance ultérieure) des utilisateurs en ligne (316*, 853). Des informations d'identification préservant la confidentialité pour identifier une personne authentique et unique en ligne pourraient minimiser ces effets indésirables sur la confidentialité (853).

Lacunes en matière de preuves : des recherches supplémentaires sont nécessaires pour étudier comment et quand l'IA à usage général risque de révéler des informations sensibles, comment l'IA à usage général peut être exécutée avec des garanties de sécurité plus solides et comment empêcher l'IA à usage général d'être utilisée pour des cas d'utilisation exploitant la vie privée. L'étendue complète des données personnelles dans les données d'entraînement de l'IA à usage général (1325) et la probabilité qu'elles soient mémorisées et exposées (831, 1364) sont inconnues et nécessitent des recherches supplémentaires. Même lorsque des données sensibles ne sont utilisées qu'au moment de l'exécution (souvent appelé « apprentissage en contexte »), des recherches supplémentaires sont nécessaires pour établir les risques de fuite d'informations par les modèles dans leur sortie (847, 1365). Lors de l'utilisation de ces systèmes d'IA à usage général, des approches cryptographiques solides pour les exécuter pourraient permettre une plus grande confidentialité et une plus grande vérifiabilité (1366), mais des travaux supplémentaires sont nécessaires pour étendre ces techniques aux grands systèmes d'IA. Pour éviter les dommages causés par des acteurs malveillants utilisant l'IA à usage général pour violer la vie privée d'autrui, des recherches supplémentaires seront nécessaires pour rendre plus difficile l'utilisation de l'IA à usage général à des fins malveillantes. De nombreuses questions techniques restent ouvertes sur la manière de préserver la confidentialité des créateurs de données, des utilisateurs et des déployeurs de systèmes d'IA tout en exploitant et en régissant l'IA à usage général (1177). De nouveaux risques pour la confidentialité peuvent également apparaître à mesure que de nouvelles capacités d'IA à usage général émergent (voir [1.3. Capacités dans les années à venir](#)).

Pour les décideurs politiques travaillant sur la protection de la vie privée, les principaux défis découlent d'un environnement technique dans lequel les méthodes visant à gérer les risques liés à la vie privée et à minimiser les dommages évoluent rapidement, dans de nombreux domaines de la chaîne d'approvisionnement de l'IA à usage général. Les domaines de risque abordés dans cette [section](#) et dans la [section 2.3.5. Les risques pour la vie privée](#) couvrent un large éventail de participants à l'écosystème de l'IA à usage général, et les stratégies d'atténuation varient en termes de faisabilité technique et de complexité (résumées dans la figure 3.3). Chaque stratégie d'atténuation imposera des coûts aux développeurs et aux déployeurs d'IA à usage général (par exemple, le nettoyage des données à l'échelle du Web est coûteux) et peut dégrader l'expérience utilisateur (par exemple, des garanties cryptographiques solides peuvent ralentir l'exécution).

(vitesses de l'IA à usage général). Ce domaine de recherche évolue et il sera difficile de prévoir dans quelle mesure les risques spécifiques liés à la confidentialité auront des stratégies d'atténuation robustes qui peuvent être déployées à grande échelle, ce qui est rendu encore plus difficile par les différences entre les communautés de l'IA et des politiques de confidentialité (822).

Méthodes pratiques pour protéger la vie privée

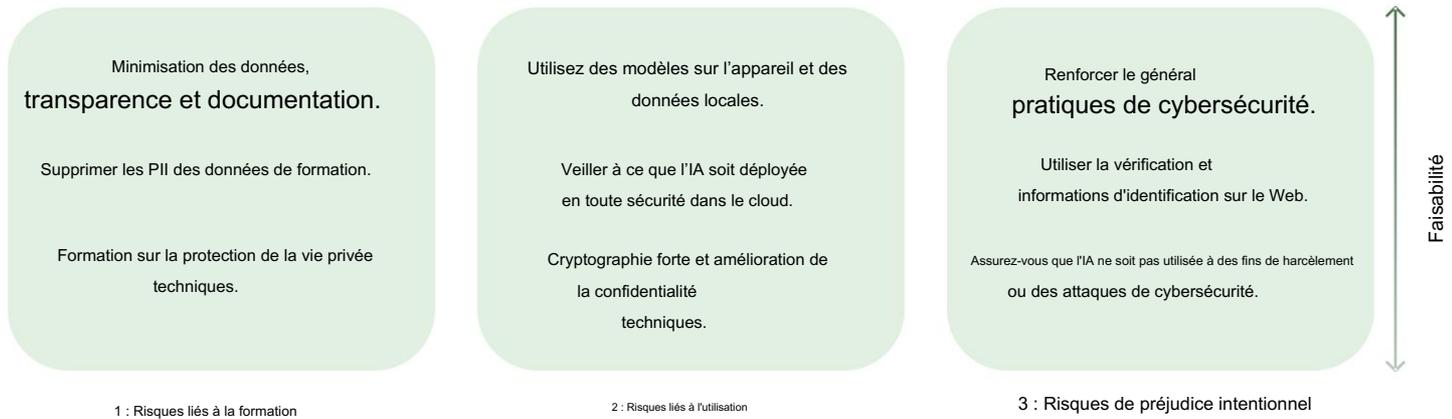


Figure 3.3 : Il existe des méthodes concrètes pour atténuer les atteintes à la vie privée causées par les systèmes d'IA à usage général, notamment la suppression des informations personnelles identifiables (IPI) des données de formation, l'utilisation de modèles intégrés aux appareils et le renforcement de la cybersécurité. Les méthodes sont classées en fonction de leur faisabilité relative au sein de chaque groupe de risques et ne sont pas exhaustives. Il existe de nombreuses mesures de protection de la vie privée et d'atténuation des atteintes, chacune présentant des niveaux de complexité et des défis de déploiement différents. Source : International AI Safety Report.

Conclusion

Le premier rapport international sur la sécurité de l'IA révèle que l'avenir de l'IA à usage général est extrêmement incertain. Il existe un large éventail de résultats possibles, même dans un avenir proche, y compris des résultats très positifs et très négatifs, ainsi que tout ce qui se situe entre les deux. L'IA à usage général présente un immense potentiel pour l'éducation, les applications médicales, les avancées de la recherche dans des domaines tels que la chimie, la biologie ou la physique, et une prospérité générale accrue grâce à l'innovation rendue possible par l'IA. S'ils sont gérés correctement, les systèmes d'IA à usage général pourraient améliorer considérablement la vie des gens dans le monde entier.

Pour tirer profit de cette technologie transformatrice en toute sécurité, les chercheurs et les décideurs politiques doivent identifier les risques qui l'accompagnent et prendre des mesures éclairées pour les atténuer. L'IA à usage général cause déjà des dommages aujourd'hui en raison d'une utilisation malveillante et de dysfonctionnements, par exemple par le biais de deepfakes, d'escroqueries et de résultats biaisés. En fonction du rythme de progression des futures capacités d'IA à usage général, des méthodes techniques que les développeurs et les régulateurs emploient pour atténuer les risques, des décisions des gouvernements et des sociétés concernant l'IA à usage général et du degré de réussite de la coordination mondiale, il est également possible que d'autres risques apparaissent. Les pires scénarios pourraient voir l'émergence de risques tels que le chômage à grande échelle, le terrorisme rendu possible par l'IA à usage général ou la perte de contrôle de l'humanité sur les systèmes d'IA à usage général. Les experts diffèrent quant à la probabilité ou à l'imminence de ces risques et à la façon dont ils interprètent les preuves existantes : certains pensent que ces risques ne se produiront que dans des décennies, tandis que d'autres pensent que l'IA à usage général pourrait entraîner de graves dangers pour la sécurité publique dans les prochaines années.

Il existe des méthodes techniques pour gérer les risques liés à l'IA à usage général, mais elles ont toutes des limites. Par exemple, les chercheurs ont mis au point des méthodes pour réduire les biais, améliorer notre compréhension du fonctionnement interne de l'IA, évaluer les capacités et les risques et rendre l'IA moins susceptible de répondre aux demandes des utilisateurs qui pourraient causer des dommages. Cependant, plusieurs caractéristiques de l'IA à usage général rendent la gestion des risques difficile. Malgré les progrès rapides des capacités, les chercheurs ne sont actuellement pas en mesure de générer des comptes-rendus compréhensibles par l'homme sur la manière dont l'IA à usage général parvient à des résultats et à des décisions. Il est donc difficile d'évaluer ou de prédire ce dont l'IA à usage général est capable et sa fiabilité, ou d'obtenir des garanties sur les risques qu'elle pourrait présenter. Les experts s'accordent largement à dire qu'il devrait être prioritaire d'améliorer notre compréhension de la manière dont l'IA à usage général parvient à des résultats et à des décisions.

L'IA ne nous arrive pas, ce sont les choix des individus qui déterminent son avenir. Comment l'IA à usage général est-elle développée et par qui, quels problèmes elle est censée résoudre, si nous serons en mesure de tirer pleinement parti de son potentiel économique, qui en bénéficie et à quels types de risques nous nous exposons – les réponses à ces questions et à bien d'autres dépendent des choix que les sociétés et les gouvernements font aujourd'hui et à l'avenir pour façonner le développement de l'IA à usage général. Étant donné que l'impact de l'IA à usage général sur de nombreux aspects de notre vie est susceptible d'être profond et que les progrès pourraient continuer à être rapides, il est urgent de travailler à un accord international et de mettre en place des mesures pour y parvenir.

Des ressources sont nécessaires pour comprendre et gérer les risques liés à cette technologie. Un débat scientifique et public constructif sera essentiel pour que les sociétés et les décideurs politiques puissent faire les bons choix.

Pour la première fois dans l'histoire, ce rapport et le rapport intermédiaire (mai 2024) ont réuni des représentants d'experts désignés par 30 pays, l'OCDE, l'UE et l'ONU, ainsi que plusieurs autres experts de renommée mondiale, afin de fournir une base scientifique commune et fondée sur des preuves pour ces discussions vitales. Nous continuons à être en désaccord sur plusieurs questions, mineures et majeures, concernant l'IA à usage général et ses capacités, ses risques et les mesures d'atténuation des risques. Cependant, nous considérons que ce rapport est essentiel pour améliorer notre compréhension collective de l'IA à usage général et de ses risques potentiels, et pour nous rapprocher d'un consensus et d'une atténuation efficace des risques, afin de garantir que l'humanité puisse profiter des avantages de l'IA à usage général en toute sécurité. Les enjeux sont élevés. Nous sommes impatients de poursuivre cet effort.

Liste des acronymes

AAVE : anglais vernaculaire afro-américain	GPQA : Évaluation de la qualité des écoles primaires
IA : intelligence artificielle	GPT : transformateur pré-entraîné génératif
OBJECTIF : Surveillance des incidents liés à l'IA	GPU : unité de traitement graphique
AIME : Concours américain de mathématiques sur invitation	HAZOP : étude de danger et d'opérabilité
Examen	PEI : pays à revenu élevé
AISI : Institut de sécurité de l'IA	TIC : information et communication
ALARP : aussi bas que raisonnablement possible	technologie
AMLAS : Assurance de l'apprentissage automatique pour une utilisation dans les systèmes autonomes	AIE : Agence internationale de l'énergie
API : interface de programmation d'application	CEI : Autorité électrotechnique internationale
ASEAN : Association des pays d'Asie du Sud-Est	Commission
Nations	OMI : Olympiade Internationale de Mathématiques
AWS : Amazon Web Services	ISO : Organisation internationale pour les Standardisation
BTWC : armes biologiques et toxiques	kWh : kilowattheure
Convention	LLM : grand modèle de langage
CBRN : chimique, biologique, radiologique et nucléaire	PRFI : pays à revenu faible et intermédiaire
CNI : infrastructures nationales critiques	MMLU : langage multitâche massif
COVID-19 : maladie à coronavirus 2019	Compréhension
CSAM : matériel d'abus sexuel sur mineur	MNIST : Institut national de l'information modifié
CTF : Capture du drapeau	Normes et technologies (base de données)
CTM : mécanisme de transition vers le charbon	MW : mégawatt
CAC : Convention sur les armes chimiques	NCII : images intimes non consensuelles
DARPA : Projet de recherche avancée sur la défense	NIST : Institut national des normes et
Agence	Technologie
DBTL : concevoir-construire-tester-apprendre	OCDE : Organisation de coopération et de développement économiques
DSIT : Département des sciences, de l'innovation et	Coopération et développement
Technologie	Ofcom : Bureau des communications
UE : Union européenne	OSS : logiciel open source
ExIF : format de fichier image échangeable	PaLM-E : Modèle de langage Pathways
FAccT : (conférence sur) l'équité,	(Incarné)
Responsabilité et transparence	PC : ordinateur personnel
FLOP : opérations en virgule flottante	PhD : Docteur en philosophie
PIB : produit intérieur brut	PII : informations personnelles identifiables
RGPD : Règlement général sur la protection des données	PPA : contrat d'achat d'électricité
GLUE : Compréhension générale du langage	ESPT : trouble de stress post-traumatique
Évaluation	PUE : efficacité énergétique
RNB : Revenu national brut	Q&R : questions et réponses
GES : gaz à effet de serre	R&D : recherche et développement
	RAG : Génération augmentée par récupération

REC : crédit énergie renouvelable

RLHF : apprentissage par renforcement à partir du feedback humain

RoG : raisonnement sur les graphes

RT : transformateur robotique

SARS-CoV-2 : coronavirus 2 du syndrome respiratoire aigu sévère

PME : petites et moyennes entreprises

SMR : petit réacteur modulaire

SOTIF : sécurité de la fonction prévue

SQLite : langage de requête structuré simplifié

SQuAD : ensemble de données de réponses aux questions de Stanford

STEM : sciences, technologie, ingénierie et mathématiques

SWE-bench : référence en ingénierie logicielle

tCO_{2e} : tonnes d'équivalent dioxyde de carbone

TPU : unité de traitement des tenseurs

TSMC : fabrication de semi-conducteurs à Taiwan

Entreprise

TWh : térawattheure

Royaume-Uni : Royaume-Uni

UNESCO : Organisation des Nations Unies pour l'éducation, la science et la culture.

Organisation scientifique et culturelle

États-Unis : États-Unis

USB : bus série universel

VD : découverte de vulnérabilité

V-JEPA : architecture prédictive d'intégration vidéo conjointe

XAI : intelligence artificielle explicable

Glossaire

Les explications ci-dessous font toutes référence à l'utilisation d'un terme relatif à l'IA.

Entraînement contradictoire : technique d'apprentissage automatique utilisée pour rendre les modèles plus fiables. Tout d'abord, les développeurs construisent des « entrées contradictoires » (par exemple via le red-teaming) qui sont conçues pour faire échouer un modèle, et ensuite, ils entraînent le modèle à reconnaître et à gérer ce type d'entrées.

Agent IA : une IA polyvalente qui peut élaborer des plans pour atteindre des objectifs, effectuer de manière adaptative des tâches impliquant plusieurs étapes et des résultats incertains en cours de route, et interagir avec son environnement. par exemple en créant des fichiers, en effectuant des actions sur le Web ou en déléguant des tâches à d'autres agents – avec peu ou pas de surveillance humaine.

Fossé en matière de recherche et développement en IA : disparité dans la recherche et le développement en IA entre les différentes régions géographiques, causée par divers facteurs, notamment une répartition inégale de la puissance de calcul, des talents, des ressources financières et des infrastructures.

Contenu factice généré par l'IA : contenu audio, textuel ou visuel, produit par l'IA générative, qui représente des personnes ou des événements d'une manière différente de la réalité de manière malveillante ou trompeuse, par exemple en montrant des personnes faisant des choses qu'elles n'ont pas faites, en disant des choses qu'elles n'ont pas dites, en changeant le lieu d'événements réels ou en représentant des événements qui ne se sont pas produits.

Cycle de vie de l'IA : les différentes étapes du développement de l'IA, y compris la collecte et le prétraitement des données, la préformation, le réglage fin, l'intégration du modèle, le déploiement, la surveillance post-déploiement et les modifications en aval.

Algorithme : un ensemble de règles ou d'instructions permettant à un système d'IA de traiter des données et d'effectuer des tâches spécifiques.

Efficacité algorithmique (d'entraînement) : ensemble de mesures de l'efficacité avec laquelle un algorithme utilise les ressources de calcul pour apprendre à partir des données, comme la quantité de mémoire utilisée ou le temps nécessaire à l'entraînement.

Transparence algorithmique : degré auquel les facteurs qui influencent les résultats de l'IA à usage général, par exemple les recommandations ou les décisions, sont connus par diverses parties prenantes. Ces facteurs peuvent inclure le fonctionnement interne du modèle d'IA, la manière dont il a été formé, les données sur lesquelles il est formé, les caractéristiques des entrées qui ont affecté ses résultats et les décisions qu'il aurait prises dans des circonstances différentes.

Alignement : propension d'une IA à utiliser ses capacités conformément aux intentions ou aux valeurs humaines. Selon le contexte, cela peut faire référence aux intentions et aux valeurs des développeurs, des opérateurs, des utilisateurs, des communautés spécifiques ou de la société dans son ensemble.

Interface de programmation d'application (API) : un ensemble de règles et de protocoles qui permet l'intégration et la communication entre les systèmes d'IA et d'autres applications logicielles.

Intelligence artificielle générale (IAG) : IA potentielle du futur qui égale ou surpasse les performances humaines sur toutes ou presque toutes les tâches cognitives.

Intelligence artificielle (IA) : domaine de l'informatique axé sur la création de systèmes ou de machines capables d'effectuer des tâches qui nécessitent généralement une intelligence humaine. Ces tâches comprennent l'apprentissage, le raisonnement, la résolution de problèmes, le traitement du langage naturel et la prise de décision.

Audit : examen formel de la conformité d'une organisation aux normes, politiques et procédures, généralement effectué par un tiers indépendant.

Automatisation : Utilisation de la technologie pour effectuer des tâches avec une intervention humaine réduite ou nulle.

Benchmark : un test ou une mesure standardisé, souvent quantitatif, utilisé pour évaluer et comparer les performances des systèmes d'IA sur un ensemble fixe de tâches conçues pour représenter une utilisation dans le monde réel.

Biais : erreurs systématiques dans les systèmes algorithmiques qui favorisent certains groupes ou visions du monde et créent souvent des résultats injustes pour certaines personnes. Les biais peuvent avoir de multiples sources, notamment des erreurs de conception algorithmique, des ensembles de données non représentatifs ou autrement erronés, ou des inégalités sociales préexistantes.

Biosécurité : ensemble de politiques, de pratiques et de mesures (par exemple, diagnostics et vaccins) conçues pour protéger les humains, les animaux, les plantes et les écosystèmes contre les agents biologiques nocifs, qu'ils soient d'origine naturelle ou introduits intentionnellement.

Capacités : l'éventail des tâches ou des fonctions qu'un système d'IA peut exécuter, et la compétence avec laquelle il peut les exécuter.

Intensité carbone : Quantité d'émissions de GES produites par unité d'énergie. Utilisée pour quantifier les émissions relatives de différentes sources d'énergie.

Compensation carbone : Compenser les émissions de GES d'une source en investissant dans d'autres activités qui empêchent des quantités comparables d'émissions ou éliminent le carbone de l'atmosphère, comme l'expansion des forêts.

Chaîne de pensée : processus de raisonnement dans lequel une IA génère des étapes intermédiaires ou des explications tout en résolvant un problème ou en répondant à une question. Cette approche imite le raisonnement logique et la délibération interne de l'homme, aidant le modèle à décomposer des tâches complexes en étapes séquentielles plus petites pour améliorer la précision et la transparence de ses résultats.

Cloud computing : paradigme de fourniture de services informatiques (notamment serveurs, stockage de données, logiciels et analyses) via Internet. Les utilisateurs peuvent accéder à ces ressources à la demande et sans infrastructure locale pour développer, former, déployer et gérer des applications d'IA.

Tâches cognitives : activités qui impliquent le traitement de l'information, la résolution de problèmes, la prise de décision et la pensée créative. Exemples : recherche, rédaction et programmation.

Calcul : abréviation de « ressources informatiques », qui fait référence au matériel (par exemple les GPU), aux logiciels (par exemple les logiciels de gestion des données) et à l'infrastructure (par exemple les centres de données) nécessaires pour former et exécuter les systèmes d'IA.

Contrôle : La capacité d'exercer une surveillance sur un système d'IA et d'ajuster ou d'arrêter son comportement s'il agit de manière indésirable.

Capacités de sape du contrôle : capacités qui, si elles étaient utilisées, permettraient à un système d'IA de saper le contrôle humain.

Droit d'auteur : une forme de protection juridique accordée aux créateurs d'œuvres originales, leur donnant le droit exclusif d'utiliser, de reproduire et de distribuer leur œuvre.

Défis CTF (Capture the Flag) : exercices souvent utilisés dans la formation en cybersécurité, conçus pour tester et améliorer les compétences des participants en les mettant au défi de résoudre des problèmes liés à la cybersécurité, comme trouver des informations cachées ou contourner les défenses de sécurité.

Centre de données : un vaste ensemble de serveurs informatiques en réseau à haute puissance utilisés pour le calcul à distance. Les centres de données hyperscale contiennent généralement plus de 5 000 serveurs.

Collecte et prétraitement des données : étape du développement de l'IA au cours de laquelle les développeurs et les travailleurs des données collectent, nettoient, étiquettent, standardisent et transforment les données de formation brutes dans un format à partir duquel le modèle peut apprendre efficacement.

Minimisation des données : pratique consistant à collecter et à conserver uniquement les données directement nécessaires à un objectif spécifique, et à les supprimer une fois cet objectif atteint.

Alignement trompeur : Désalignement difficile à détecter, car le système se comporte d'une manière qui, au moins au départ, semble bénigne.

Deepfake : un type de faux contenu généré par l'IA, composé de contenu audio ou visuel, qui déforme les faits et gestes de personnes réelles en les faisant croire qu'elles n'ont pas réellement fait ou dit quelque chose.

Apprentissage profond : technique d'apprentissage automatique dans laquelle de grandes quantités de données et de calculs sont utilisées pour former des réseaux neuronaux artificiels multicouches (inspirés du cerveau biologique) afin d'apprendre et d'extraire automatiquement des fonctionnalités de haut niveau à partir de grands ensembles de données, permettant ainsi de puissantes capacités de reconnaissance de modèles et de prise de décision.

Défense en profondeur : stratégie qui comprend la superposition de plusieurs mesures d'atténuation des risques dans les cas où aucune méthode existante ne peut assurer la sécurité.

Déploiement : processus de mise en œuvre de systèmes d'IA dans des applications, des produits ou des services du monde réel où ils peuvent répondre à des demandes et fonctionner dans un contexte plus large.

Développeur : Toute organisation qui conçoit, construit, intègre, adapte ou combine des modèles ou des systèmes d'IA.

Fossé numérique : disparité dans l'accès aux technologies de l'information et de la communication (TIC), notamment à Internet, entre différentes régions géographiques ou groupes de personnes.

Informatique légale : processus de traçage de l'origine et de la diffusion des médias numériques.

Infrastructure numérique : les services et installations fondamentaux nécessaires au fonctionnement des technologies numériques, notamment le matériel, les logiciels, les réseaux, les centres de données et les systèmes de communication.

Discrimination : Traitement injuste d'individus ou de groupes en raison de leurs attributs, tels que la race, le sexe, l'âge, la religion ou d'autres caractéristiques protégées.

Désinformation : information fautive ou trompeuse générée ou diffusée dans le but de tromper ou d'influencer des personnes. Voir « Désinformation » pour plus de contraste.

Formation distribuée : processus de formation de modèles d'IA sur plusieurs processeurs et serveurs, concentrés dans un ou plusieurs centres de données.

Science à double usage : recherche et technologie qui peuvent être appliquées à des fins bénéfiques, telles que la médecine ou les solutions environnementales, mais qui peuvent aussi être potentiellement utilisées à mauvais escient pour causer des dommages, comme dans le développement d'armes biologiques ou chimiques.

Comportement émergent : capacité des systèmes d'IA à agir d'une manière qui n'a pas été explicitement programmée ou prévue par leurs développeurs ou utilisateurs.

Évaluations : Évaluations systématiques des performances, des capacités, des vulnérabilités ou des impacts potentiels d'un système d'IA. Les évaluations peuvent inclure des analyses comparatives, des red-teams et des audits et peuvent être menées avant et après le déploiement du modèle.

IA explicable (XAI) : programme de recherche visant à créer des systèmes d'IA qui fournissent des explications claires et compréhensibles de leurs décisions, permettant aux utilisateurs de comprendre comment et pourquoi des résultats spécifiques sont générés.

Équité : valeur sociétale selon laquelle les systèmes d'IA doivent prendre des décisions exemptes de préjugés ou de discrimination injuste, en traitant tous les individus et tous les groupes de manière équitable, notamment en ce qui concerne les attributs protégés tels que la race, le sexe, l'âge ou le statut socio-économique.

Fair use : doctrine juridique américaine qui offre une défense contre les réclamations pour violation du droit d'auteur pour des utilisations limitées de matériel protégé par le droit d'auteur sans autorisation à des fins telles que la critique, le commentaire, le reportage d'actualité, l'éducation et la recherche. Certains autres pays autorisent des droits d'utilisation similaires sous le nom de « fair handling ».

Tests sur le terrain : pratique consistant à évaluer les risques de l'IA à usage général dans des conditions réelles.

Réglage fin : processus consistant à adapter un modèle d'IA pré-entraîné à une tâche spécifique ou à le rendre plus utile en général en l'entraînant sur des données supplémentaires.

Avantage du premier entrant : l'avantage concurrentiel obtenu en étant le premier à établir une position de marché significative dans un secteur.

FLOP : « Opérations en virgule flottante » : nombre d'opérations de calcul effectuées par un programme informatique. Souvent utilisé comme mesure de la quantité de calcul utilisée pour entraîner un modèle d'IA.

Modèle de fondation : un modèle d'IA à usage général conçu pour être adaptable à un large éventail de tâches en aval.

Frontier AI : terme parfois utilisé pour désigner une IA particulièrement performante, qui égale ou dépasse les capacités de l'IA la plus avancée d'aujourd'hui. Aux fins du présent rapport, l'IA de pointe peut être considérée comme une IA polyvalente particulièrement performante.

IA à usage général : systèmes d'IA conçus pour effectuer une large gamme de tâches dans divers domaines, plutôt que d'être spécialisés dans une fonction spécifique. Voir « IA restreinte » pour plus de contraste.

IA générative : IA capable de créer du nouveau contenu tel que du texte, des images ou de l'audio en apprenant des modèles à partir de données existantes et en générant de nouveaux résultats qui reflètent ces modèles.

Émissions de GES (gaz à effet de serre) : rejet de gaz tels que le dioxyde de carbone (CO₂), le méthane, l'oxyde nitreux et les hydrofluorocarbures qui créent une barrière emprisonnant la chaleur dans l'atmosphère. Indicateur clé du changement climatique.

Travail fantôme : travail caché effectué par les travailleurs pour soutenir le développement et le déploiement de modèles ou de systèmes d'IA (par exemple via l'étiquetage des données).

Généralisation erronée des objectifs : situation dans laquelle un système d'IA suit correctement un objectif dans son environnement de formation, mais l'applique de manière inattendue lorsqu'il fonctionne dans un environnement différent.

Mauvaise spécification des objectifs : inadéquation entre l'objectif donné à une IA et l'intention du développeur, conduisant l'IA à adopter des comportements imprévus ou indésirables.

GPU (unité de traitement graphique) : une puce informatique spécialisée, conçue à l'origine pour l'infographie, qui est désormais largement utilisée pour gérer des tâches de traitement parallèle complexes essentielles à la formation et à l'exécution de modèles d'IA.

Garde-fous : contraintes de sécurité intégrées pour garantir qu'un système d'IA fonctionne comme souhaité et évite les conséquences néfastes.

Piratage informatique : acte consistant à exploiter les vulnérabilités ou les faiblesses d'un système informatique, d'un réseau ou d'un logiciel pour obtenir un accès non autorisé, manipuler des fonctionnalités ou extraire des informations.

Hallucination : Informations inexacts ou trompeuses générées par un système d'IA, par exemple des faits ou des citations erronés.

Porte dérobée matérielle : fonctionnalité d'un appareil, intentionnellement ou non, créée par un fabricant ou un tiers, qui peut être utilisée pour contourner les protections de sécurité afin de surveiller, contrôler ou extraire des données à l'insu de l'utilisateur.

Danger : Tout événement ou activité susceptible de causer un préjudice, tel qu'une perte de vie, une blessure, une perturbation sociale ou des dommages environnementaux.

Pays à revenu élevé (PRE) : Pays dont le revenu national brut (RNB) par habitant est supérieur à 14 005 dollars, tel que calculé par la Banque mondiale.

L'humain dans la boucle : une exigence selon laquelle les humains doivent superviser et approuver les processus automatisés dans des domaines critiques.

Engagements de type « si-alors » : accords conditionnels, cadres ou réglementations qui spécifient les actions ou les obligations à réaliser lorsque certaines conditions prédéfinies sont remplies.

Signalement d'incident : documenter et partager les cas dans lesquels le développement ou le déploiement de l'IA a causé des dommages directs ou indirects.

Inférence : processus par lequel une IA génère des sorties basées sur une entrée donnée, appliquant ainsi les connaissances acquises pendant la formation.

Améliorations du temps d'inférence : techniques utilisées pour améliorer les performances d'un système d'IA après son entraînement initial, sans modifier le modèle sous-jacent. Cela comprend des méthodes d'invite intelligentes, des méthodes de sélection de réponses (par exemple, l'échantillonnage de plusieurs réponses et le choix d'une réponse majoritaire), l'écriture de longues « chaînes de pensée », l'« échafaudage » des agents, etc.

Entrée (à un système d'IA) : les données ou l'invite soumises à un système d'IA, comme du texte ou une image, que le système d'IA traite et transforme en sortie.

Transparence institutionnelle : degré auquel les entreprises d'IA divulguent des informations techniques ou organisationnelles à l'examen public ou gouvernemental, y compris les données de formation, les architectures de modèles, les données d'émissions, les mesures de sécurité et de sûreté ou les processus décisionnels.

Propriété intellectuelle : Créations de l'esprit sur lesquelles des droits légaux peuvent être accordés, y compris les œuvres littéraires et artistiques, les symboles, les noms et les images.

Interprétabilité : degré auquel les humains peuvent comprendre le fonctionnement interne d'un modèle d'IA, notamment pourquoi il a généré un résultat ou une décision particulière. Un modèle est hautement interprétable si ses processus mathématiques peuvent être traduits en concepts qui permettent aux humains de retracer les facteurs et la logique spécifiques qui ont influencé le résultat du modèle.

Recherche d'interprétabilité : étude du fonctionnement interne des modèles d'IA à usage général et développement de méthodes pour rendre cela compréhensible pour les humains.

Jailbreaking : Génération et envoi d'invites conçues pour contourner les garde-fous et amener un système d'IA à produire du contenu nuisible, comme des instructions pour la construction d'armes.

Marché du travail : système dans lequel les employeurs cherchent à embaucher des travailleurs et les travailleurs cherchent un emploi, englobant la création d'emplois, la perte d'emplois et les salaires.

Perturbation du marché du travail : changements importants et souvent complexes sur le marché du travail qui affectent la disponibilité des emplois, les compétences requises, la répartition des salaires ou la nature du travail dans tous les secteurs et professions.

Modèle linguistique de grande taille (LLM) : un modèle d'IA formé sur de grandes quantités de données textuelles pour effectuer des tâches de traitement du langage, telles que la génération, la traduction ou la synthèse de texte.

Droits à l'image : droits qui protègent l'image, la voix, le nom ou d'autres aspects identifiables d'une personne contre toute utilisation commerciale non autorisée.

Scénario de perte de contrôle : scénario dans lequel un ou plusieurs systèmes d'IA à usage général se mettent à fonctionner en dehors du contrôle de quiconque, sans voie claire pour reprendre le contrôle.

Pays à revenu faible ou intermédiaire (PRFI) : pays dont le revenu national brut (RNB) par habitant est inférieur à 14 005 dollars, tel que calculé par la Banque mondiale.

Apprentissage automatique (ML) : un sous-ensemble de l'IA axé sur le développement d'algorithmes et de modèles qui apprennent à partir des données et améliorent leurs performances sur les tâches au fil du temps sans être explicitement programmés.

Dysfonctionnement : L'incapacité d'un système d'IA à usage général à fonctionner comme prévu par son développeur ou son utilisateur, entraînant des résultats incorrects ou nuisibles ou des perturbations opérationnelles.

Utilisation malveillante : utilisation de l'IA pour causer intentionnellement des dommages.

Logiciel malveillant : logiciel nuisible conçu pour endommager, perturber ou obtenir un accès non autorisé à un système informatique. Il comprend les virus, les logiciels espions et autres programmes malveillants qui peuvent voler des données ou causer des dommages.

Risque marginal : Le risque supplémentaire introduit par un modèle ou un système d'IA à usage général par rapport à une référence pertinente, comme un risque comparable posé par une technologie non IA existante.

Concentration du marché : degré auquel un petit nombre d'entreprises contrôlent un secteur, ce qui entraîne une réduction de la concurrence et un contrôle accru sur les prix et l'innovation.

Compréhension du langage multitâche massif (MMLU) : une référence largement utilisée dans la recherche en IA qui évalue les performances d'un modèle d'IA à usage général sur un large éventail de tâches et de sujets zones.

Désalignement : propension d'une IA à utiliser ses capacités d'une manière qui entre en conflit avec les intentions ou les valeurs humaines. Selon le contexte, cela peut faire référence aux intentions et aux valeurs des développeurs, des opérateurs, des utilisateurs, des communautés spécifiques ou de la société dans son ensemble.

Désinformation : information fausse ou trompeuse générée ou diffusée sans intention de tromper. Voir « Désinformation » pour plus de contraste.

Modalités : Les types de données qu'un système d'IA peut recevoir avec compétence en entrée et produire en sortie, y compris du texte (langage ou code), des images, des vidéos et des actions robotiques.

Modèle : programme informatique, souvent basé sur l'apprentissage automatique, conçu pour traiter des entrées et générer des sorties. Les modèles d'IA peuvent effectuer des tâches telles que la prédiction, la classification, la prise de décision ou la génération, constituant ainsi le cœur des applications d'IA.

Fiche modèle : document fournissant des informations utiles sur un modèle d'IA, par exemple sur son objectif, ses directives d'utilisation, ses données de formation, ses performances sur des tests de référence ou ses fonctionnalités de sécurité.

Publication du modèle : mise à disposition d'un modèle d'IA formé pour que les entités en aval puissent l'utiliser, l'étudier ou le modifier davantage, ou pour l'intégrer dans leurs propres systèmes.

IA restreinte : type d'IA spécialisé pour effectuer une tâche spécifique ou quelques tâches très similaires, telles que le classement des résultats de recherche sur le Web, la classification des espèces animales ou le jeu d'échecs. Voir « IA à usage général » pour plus de contraste.

Réseau neuronal : type de modèle d'IA constitué d'une structure mathématique inspirée du cerveau humain et composée de nœuds interconnectés (comme des neurones) qui traitent et apprennent à partir des données. Les systèmes d'IA polyvalents actuels sont basés sur des réseaux neuronaux.

Domaines ouverts : environnements dans lesquels des systèmes d'IA peuvent être déployés et qui présentent un très large éventail de scénarios possibles. Dans les domaines ouverts, les développeurs ne peuvent généralement pas anticiper et tester toutes les manières possibles dont un système d'IA pourrait être utilisé.

Modèle à pondération ouverte : un modèle d'IA dont les pondérations sont disponibles publiquement en téléchargement, comme Llama ou Stable Diffusion. Les modèles à pondération ouverte peuvent être, mais ne sont pas nécessairement, open source.

Modèle open source : un modèle d'IA publié en téléchargement public sous une licence open source.

La licence open source accorde la liberté d'utiliser, d'étudier, de modifier et de partager le modèle pour n'importe quel

Objectif. Il subsiste un certain désaccord quant aux composants du modèle (poids, code, données de formation) et à la documentation qui doivent être accessibles au public pour que le modèle soit qualifié d'open source.

Paramètres : les composants numériques d'un modèle d'IA, tels que les pondérations et les biais, qui sont appris à partir des données pendant l'entraînement et qui déterminent la manière dont le modèle traite les entrées pour générer les sorties. Notez que le terme « biais » est ici un terme mathématique qui n'est pas lié au biais dans le contexte de la discrimination.

Pathogène : Un micro-organisme, par exemple un virus, une bactérie ou un champignon, qui peut provoquer une maladie chez les humains, les animaux ou les plantes.

Tests de pénétration : pratique de sécurité dans laquelle des experts agréés ou des systèmes d'IA simulent des cyberattaques sur un système informatique, un réseau ou une application pour évaluer proactivement sa sécurité. L'objectif est d'identifier et de corriger les faiblesses avant qu'elles ne soient exploitées par de véritables attaquants.

Informations personnelles identifiables (IPI) : toute donnée permettant d'identifier directement ou indirectement une personne (par exemple, un nom ou un numéro d'identification). Comprend les informations qui peuvent être utilisées seules ou combinées à d'autres données pour identifier une personne de manière unique.

Surveillance post-déploiement : processus par lesquels les développeurs d'IA suivent l'impact du modèle et les mesures de performance, collectent et analysent les commentaires des utilisateurs et apportent des améliorations itératives pour résoudre les problèmes ou les limitations découverts lors d'une utilisation réelle.

Pré-formation : étape de développement d'un modèle d'IA à usage général au cours de laquelle les modèles apprennent des modèles à partir de grandes quantités de données. Il s'agit de l'étape de développement d'un modèle qui requiert le plus de calculs.

Confidentialité : droit d'une personne ou d'un groupe à contrôler la manière dont les autres accèdent à leurs informations et activités sensibles ou les traitent.

Invite : une entrée dans un système d'IA, telle qu'une question ou une requête textuelle, que le système traite et à laquelle il répond.

Course vers le bas : un scénario concurrentiel dans lequel des acteurs comme les entreprises ou les États-nations privilégient le développement rapide de l'IA plutôt que la sécurité.

Ransomware : un type de logiciel malveillant qui verrouille ou crypte les fichiers ou le système d'un utilisateur, les rendant inaccessibles jusqu'à ce qu'une rançon (généralement de l'argent) soit versée à l'attaquant.

Effet de rebond : En économie, réduction des améliorations attendues en raison d'une augmentation de l'efficacité, résultant de changements corrélés dans le comportement, les habitudes d'utilisation ou d'autres changements systémiques.

Par exemple, une amélioration de 25 % de l'efficacité du moteur à combustion automobile (km/litre) entraînera une réduction de moins de 25 % des émissions, car la réduction correspondante du coût de l'essence par kilomètre parcouru rendra moins cher le fait de conduire davantage, ce qui limitera les améliorations.

Red teaming : processus systématique dans lequel des individus ou des équipes dédiés recherchent des vulnérabilités, des limites ou des risques d'utilisation abusive par le biais de diverses méthodes. Souvent, l'équipe rouge recherche des entrées qui induisent un comportement indésirable dans un modèle ou un système afin d'identifier les failles de sécurité.

Apprentissage par renforcement à partir de commentaires humains (RLHF) : une technique d'apprentissage automatique dans laquelle un modèle d'IA est affiné en utilisant des évaluations ou des préférences fournies par l'homme comme signal de récompense, permettant

le système pour apprendre et ajuster son comportement pour mieux s'aligner sur les valeurs et les intentions humaines grâce à un entraînement itératif.

Fiabilité : capacité d'un système d'IA à exécuter systématiquement sa fonction prévue.

Politique de mise à l'échelle responsable (RSP) : un ensemble de protocoles techniques et organisationnels, généralement sous un format « si-alors » pour différents niveaux de capacité, qui spécifient des règles pour le développement et le déploiement sûrs de systèmes d'IA de plus en plus performants.

Génération augmentée par récupération (RAG) : une technique qui permet aux LLM d'extraire des informations d'autres sources lors de l'inférence, telles que les résultats de recherche sur le Web ou une base de données interne d'une entreprise, permettant ainsi des réponses plus précises ou personnalisées.

Risque : Combinaison de la probabilité et de la gravité d'un préjudice résultant du développement, du déploiement ou de l'utilisation de l'IA.

Facteurs de risque : propriétés ou conditions qui peuvent accroître les risques d'un système d'IA. Par exemple, des garde-fous faibles constituent un facteur de risque qui pourrait permettre à des acteurs malveillants d'utiliser un système d'IA pour une cyberattaque.

Gestion des risques : processus systématique d'identification, d'évaluation, d'atténuation et de surveillance des risques.

Seuil de risque : limite quantitative ou qualitative qui distingue les risques acceptables des risques inacceptables et déclenche des actions spécifiques de gestion des risques lorsqu'elle est dépassée.

Tolérance au risque : niveau de risque qu'un individu ou une organisation est prêt à assumer.

Robustesse (d'un système d'IA) : propriété de se comporter en toute sécurité dans un large éventail de circonstances.

Sécurité (d'un système d'IA) : La propriété d'éviter les résultats nuisibles, tels que la fourniture d'informations dangereuses aux utilisateurs, l'utilisation à des fins néfastes ou des dysfonctionnements coûteux dans des contextes à enjeux élevés.

Dossier de sécurité : argument structuré, généralement produit par un développeur et étayé par des preuves, selon lequel un modèle ou un système d'IA est suffisamment sûr dans un contexte opérationnel donné. Les développeurs ou les régulateurs peuvent utiliser les dossiers de sécurité comme base pour des décisions importantes (par exemple, pour savoir s'il faut déployer un système d'IA).

Échafaudage : logiciel supplémentaire construit autour d'un système d'IA qui l'aide à effectuer une tâche. Par exemple, un système d'IA peut avoir accès à une application de calcul externe pour augmenter ses performances sur des problèmes arithmétiques. Un échafaudage plus sophistiqué peut structurer les résultats d'un modèle et guider le modèle pour améliorer ses réponses étape par étape.

Lois d'échelle : relations systématiques observées entre la taille d'un modèle d'IA (ou la quantité de temps, de données ou de ressources informatiques utilisées dans la formation ou l'inférence) et ses performances.

Sécurité (d'un système d'IA) : propriété d'être résilient aux interférences techniques, telles que les cyberattaques ou les fuites du code source du modèle sous-jacent.

Semi-conducteur : un matériau (généralement du silicium) dont les propriétés électriques peuvent être contrôlées avec précision, constituant l'élément fondamental des puces informatiques, telles que les GPU.

Données sensibles : informations qui, si elles étaient divulguées ou mal traitées, pourraient entraîner un préjudice, un embarras, un inconfort ou une injustice pour une personne ou une organisation.

Point de défaillance unique : élément d'un système plus vaste dont la défaillance perturbe l'ensemble du système. Par exemple, si un seul système d'IA joue un rôle central dans l'économie ou une infrastructure critique, son dysfonctionnement pourrait provoquer des perturbations généralisées dans toute la société.

Données synthétiques : données telles que du texte ou des images générées artificiellement, par exemple par des systèmes d'IA à usage général. Les données synthétiques peuvent être utilisées pour former des systèmes d'IA, par exemple lorsque les données naturelles de haute qualité sont rares.

Système : une configuration intégrée qui combine un ou plusieurs modèles d'IA avec d'autres composants, tels que des interfaces utilisateur ou des filtres de contenu, pour produire une application avec laquelle les utilisateurs peuvent interagir.

Intégration système : processus consistant à combiner un modèle d'IA avec d'autres composants logiciels pour produire un « système d'IA » complet et prêt à l'emploi. Par exemple, l'intégration peut consister pour les développeurs à combiner un modèle d'IA à usage général avec des filtres de contenu et une interface utilisateur pour produire une application de chatbot.

Risques systémiques : risques sociétaux plus larges associés au développement et au déploiement d'IA à usage général, au-delà des capacités des modèles ou systèmes individuels. Les exemples de risques systémiques vont des répercussions potentielles sur le marché du travail aux atteintes à la vie privée et aux dommages environnementaux.

Il convient de noter que cela diffère de la définition du « risque systémique » donnée par la loi sur l'IA de l'Union européenne.

Là, le terme fait référence à « un risque spécifique aux capacités à fort impact des modèles d'IA à usage général, ayant un impact significatif ».

Toxine : Substance toxique produite par des organismes vivants (tels que des bactéries, des plantes ou des animaux), ou créée synthétiquement pour imiter une toxine naturelle, qui peut provoquer des maladies, des dommages ou la mort chez d'autres organismes en fonction de sa puissance et du niveau d'exposition.

TPU (tensor processing unit) : une puce informatique spécialisée, développée par Google pour accélérer les charges de travail d'apprentissage automatique, qui est désormais largement utilisée pour gérer des calculs à grande échelle pour la formation et l'exécution de modèles d'IA.

Marque déposée : Un symbole, un mot ou une phrase légalement enregistré ou établi par l'usage pour représenter une entreprise ou un produit, le distinguant des autres sur le marché.

Transformer : une architecture de modèle d'apprentissage profond (réseau neuronal) au cœur de la plupart des modèles d'IA modernes à usage général. L'architecture du transformateur s'est avérée particulièrement efficace pour convertir des volumes de plus en plus importants de données d'apprentissage et de puissance de calcul en de meilleures performances de modèle.

Filigrane : un motif subtil, souvent imperceptible, intégré dans un contenu généré par l'IA (tel que du texte, des images ou de l'audio) pour indiquer son origine artificielle, vérifier sa source ou détecter une éventuelle utilisation abusive.

Exploration Web : utilisation d'un programme automatisé, souvent appelé robot d'exploration ou bot, pour naviguer sur le Web, dans le but de collecter des données à partir de sites Web.

Pondérations : paramètres de modèle qui représentent la force de connexion entre les nœuds d'un réseau neuronal. Les pondérations jouent un rôle important dans la détermination de la sortie d'un modèle en réponse à une entrée donnée et sont mises à jour de manière itérative pendant l'entraînement du modèle pour améliorer ses performances.

Dénonciation : divulgation d'informations, par un membre individuel d'une organisation, sur des activités illégales ou contraires à l'éthique se déroulant au sein de l'organisation, à des autorités internes ou externes ou au public.

Le gagnant rafle tout : Concept économique faisant référence aux cas dans lesquels une seule entreprise capte une très grande part de marché, même si les consommateurs ne préfèrent que légèrement ses produits ou services à ceux de ses concurrents.

Vulnérabilité zero-day : faille de sécurité non découverte ou non corrigée dans un logiciel ou un matériel. Les attaquants pouvant déjà l'exploiter, les développeurs disposent d'un délai de « zéro jour » pour la corriger.

Comment citer ce rapport

Citation formatée

Y. Bengio, S. Mindermann, D. Privitera, T. Besiroglu, R. Bommasani, S. Casper, Y. Choi, P. Fox, B. Garfinkel, D. Goldfarb, H. Heidari, A. Ho, S. Kapoor, L. Khalatbari, S. Longpre, S. Manning, V. Mavroudis, M. Mazeika, J. Michael, J. Newman, KY Ng, CT Okolo, D. Raji, G. Sastry, E. Seger, T. Skeadas, T. South, E. Strubell, F. Tramèr, L. Velasco, N. Wheeler, D. Acemoglu, O. Adekanmbi, D. Dalrymple, TG Dietterich, P. Fung, P.-O. Gourinchas, F. Heintz, G. Hinton, N. Jennings, A. Krause, S. Leavy, P. Liang, T. Ludermit, V. Marda, H. Margetts, J. McDermid, J. Munga, A. Narayanan, A. Nelson, C. Neppel, A. Oh, G. Ramchurn, S. Russell, M. Schaake, B. Schölkopf, D. Song, A. Soto, L. Tiedrich, G. Varoquaux, EW Felten, A. Yao, Y.-Q. Zhang, O. Ajala, F. Albalawi, M. Alserkal, G. Avrin, C. Busch, AC P. de LF de Carvalho, B. Fox, AS Gill, AH Hatip, J. Heikkilä, C. Johnson, G. Jolly, Z. Katzir, SM Khan, H. Kitano, A. Krüger, KM Lee, DV Ligt, JR López Portillo, D., O. Molchanovskiy, A. Monti, N. Mwamanzi, M. Nemer, N. Oliver, R. Pezoa Rivera, B. Ravindran, H. Riza, C. Rugege, C. Seoighe, H. Sheikh, J. Sheehan, D. Wong, Y. Zeng, « Rapport international sur la sécurité de l'IA » (DSIT 2025/001, 2025); <https://www.gov.uk/government/publications/international-ai-safety-report-2025>

Entrée Bibtex

@techreport{ISRSAA2025,

titre = {Rapport international sur la sécurité de l'IA},

auteur = {Bengio, Yoshua et Mindermann, Soren et Privitera, Daniel et Besiroglu, Tamay et Bommasani, Rishi et Casper, Stephen et Choi, Yejin et Fox, Philip et Garfinkel, Ben et Goldfarb, Danielle et Heidari, Hoda et Ho, Anson et Kapoor, Sayash et Khalatbari, Leila et Longpre, Shayne et Manning, Sam et Mavroudis, Vasilios et Mazeika, Mantas et Michael, Julian et Newman, Jessica et Ng, Kwan Yee et Okolo, Chinasa T. et Raji, Deborah et Sastry, Girish et Seger, Elizabeth et Skeadas, Theodora et South, Tobin et Strubell, Emma et Tramèr, Florian et Velasco, Lucia et Wheeler, Nicole et Acemoglu, Daron et Adekanmbi, Olubayo et Dalrymple, David et Dietterich, Thomas G. et Felten, Edward W. et Fung, Pascale et Gourinchas, Pierre-Olivier et Heintz, Fredrik et Hinton, Geoffrey et Jennings, Nick et Krause, Andreas et Leavy, Susan et Liang, Percy et Ludermit, Teresa et Marda, Vidushi et Margetts, Helen et McDermid, John et Munga, Jane et Narayanan, Arvind et Nelson, Alondra et Neppel, Clara et Oh, Alice et Ramchurn, Gopal et Russell, Stuart et Schaake, Marietje et Schölkopf, Bernhard et Song, Dawn et Soto, Alvaro et Tiedrich, Lee et Varoquaux, Ga'e et Yao, Andrew et Zhang, Ya-Qin et Ajala, Olubunmi et Albalawi, Fahad et Alserkal, Marwan et Avrin, Guillaume et Busch, Christian et {de Carvalho}, Andr'e Carlos Ponce de Leon Ferreira et Fox, Bronwyn et Gill, Amandeep Singh et Hatip, Ahmet Halit et Heikkil'a}, Juha et Johnson, Chris et Jolly, Gill et Katzir, Ziv et Khan, Saif M. et Kitano, Hiroaki et Kr'u ger, Antonio et Lee, Kyoung Mu et Ligt, Dominic Vincent et {L'opez Portillo}, Jos'e Ram'on et Molchanovskiy, Oleksii et Monti, Andrea et Mwamanzi, Nusu et Nemer, Mona et Oliver, Nuria et {Pezoa Rivera}, Raquel et Ravindran, Balaraman et Riza, Hammam et Rugege, Crystal et Seoighe, Ciar'a'n et Sheehan, Jerry et Sheikh, Haroon et Wong, Denise et Zeng, Yi},

année = {2025},

numéro = {DSIT 2025/001},

URL = {<https://www.gov.uk/government/publications/rapport-international-sur-la-securite-de-l-IA-2025>}

}

Références

* Indique que la référence était un rapport publié par une société d'IA à but lucratif ou qu'au moins 50 % des auteurs d'une prépublication (en fonction de leurs affiliations répertoriées) travaillent pour une société d'IA à but lucratif. Cette classification est basée uniquement sur les données d'affiliation fournies dans les publications, est fournie à titre informatif uniquement et ne doit pas être considérée comme exhaustive.

- 1 R. Simmons-Edler, R. Badman, S. Longpre, K. Rajan, « Les armes autonomes alimentées par l'IA risquent d'entraîner une instabilité géopolitique et de menacer la recherche sur l'IA » dans Actes de la 41e Conférence internationale sur l'apprentissage automatique (ICML 2024) (PMLR, 2024) ; <https://proceedings.mlr.press/v235/simmons-edler24a.html>.
- 2* OpenAI, « Carte système OpenAI o1 » (OpenAI, 2024) ; <https://cdn.openai.com/o1-system-card-20240917.pdf>.
- 3* OpenAI, « Carte système GPT-4o » (OpenAI, 2024) ; <https://cdn.openai.com/gpt-4o-system-card.pdf>.
- 4* Equipe Gêmeaux, R. Anil, S. Borgeaud, J.-B. Alayrac, J. Yu, R. Soricut, J. Schalkwyk, AM Dai, A. Hauth, K. Millican, D. Silver, M. Johnson, I. Antonoglou, J. Schrittwieser, A. Glaese, J. Chen, E. Pitler, ... O. Vinyals, « Gemini : une famille de modèles multimodaux hautement performants » (Google DeepMind, 2023) ; <http://arxiv.org/abs/2312.11805>.
- 5* Anthropic, Claude 3.5 Sonnet Modèle de carte Addendum (2024) ; https://www-cdn.anthropic.com/fed9cc193a14b84131812372d8d5857f8f304c52/Model_Card_Claude_3_Addendum.pdf.
- 6* Cohere, Commande R+ (2024) ; <https://docs.cohere.com/v2/docs/command-r-plus>.
- 7* B. Hui, J. Yang, Z. Cui, J. Yang, D. Liu, L. Zhang, T. Liu, J. Zhang, B. Yu, K. Lu, K. Dang, Y. Fan, Y. Zhang, A. Yang, R. Men, F. Huang, B. Zheng, ... J. Lin, rapport technique Qwen2.5-Coder, arXiv [cs.CL] (2024) ; <http://arxiv.org/abs/2409.12186>.
- 8* Y. Sun, S. Wang, S. Feng, S. Ding, C. Pang, J. Shang, J. Liu, X. Chen, Y. Zhao, Y. Lu, W. Liu, Z. Wu, W. Gong, J. Liang, Z. Shang, P. Sun, W. Liu, ... H. Wang, ERNIE 3.0 : pré-formation à grande échelle améliorée par les connaissances pour la compréhension et la génération de langage, arXiv [cs.CL] (2021) ; <http://arxiv.org/abs/2107.02137>.
- 9* X. Sun, Y. Chen, Y. Huang, R. Xie, J. Zhu, K. Zhang, S. Li, Z. Yang, J. Han, X. Shu, J. Bu, Z. Chen, X. Huang, F. Lian, S. Yang, J. Yan, Y. Zeng, ... J. Jiang, Hunyuan-Large : Un modèle MoE open source avec 52 milliards de paramètres activés par Tencent, arXiv [cs.CL] (2024) ; <http://arxiv.org/abs/2411.02265>.
- 10* 01.AI, A. Young, B. Chen, C. Li, C. Huang, G. Zhang, G. Zhang, H. Li, J. Zhu, J. Chen, J. Chang, K. Yu, P. Liu, Q. Liu, S. Yue, S. Yang, S. Yang, ... Z. Dai, Yi : Modèles de fondation ouverts par 01.AI, arXiv [cs.CL] (2024) ; <http://arxiv.org/abs/2403.04652>.
- 11* Meta, Llama-3.1-8B Carte modèle officielle (2024) ; <https://huggingface.co/meta-llama/Llama-3.1-8B>.
- 12* Mistral AI, carte modèle pour Mistral-Large-Instruct-2407 (2024) ; <https://huggingface.co/mistralai/Mistral-Large-Instruct-2407>.
- 13 L. Yang, Z. Zhang, Y. Song, S. Hong, R. Xu, Y. Zhao, W. Zhang, B. Cui, M.-H. Yang, Modèles de diffusion : une étude approfondie des méthodes et des applications. ACM Computing Surveys 56, 1–39 (2023) ; <https://doi.org/10.1145/3626235>.
- 14* OpenAI, « Carte système DALL-E 3 » (OpenAI, 2023) ; https://cdn.openai.com/papers/DALL_E_3_System_Card.pdf.
- 15* P. Esser, S. Kulal, A. Blattmann, R. Entezari, J. Müller, H. Saini, Y. Levi, D. Lorenz, A. Sauer, F. Boesel, D. Podell, T. Dockhorn, Z. English, K. Lacey, A. Goodwin, Y. Marek, R. Rombach, Mise à l'échelle des transformateurs de flux rectifiés pour la synthèse d'images haute résolution, arXiv [cs.CV] (2024) ; <http://arxiv.org/abs/2403.03206>.
- 16* T. Brooks, B. Peebles, C. Holmes, W. DePue, Y. Guo, L. Jing, D. Schnurr, J. Taylor, T. Luhman, E. Luhman, C. Ng, R. Wang, A. Ramesh, « Modèles de génération vidéo comme simulateurs du monde » (OpenAI, 2024) ; <https://openai.com/research/video-generation-models-as-world-simulators>.
- 17 B. Guo, X. Shan, J. Chung, Une étude comparative sur les caractéristiques et les applications des outils d'IA - Focus sur PIKA Labs et RUNWAY. Revue internationale de l'Internet, de la radiodiffusion et de la communication 16, 86–91 (2024) ; <https://doi.org/10.7236/ijbc.2024.16.1.86>.
- 18 D. Driess, F. Xia, MSM Sajjadi, C. Lynch, A. Chowdhery, B. Ichter, A. Wahid, J. Tompson, Q. Vuong, T. Yu, W. Huang, Y. Chebotar, P. Sermanet, D. Duckworth, S. Levine, V. Vanhoucke, K. Hausman, ... P. Florence, « PaLM-E : un modèle de langage multimodal incarné » dans Actes de la 40e Conférence internationale sur l'apprentissage automatique (ICML'23) (PMLR, Honolulu, Hawaï, États-Unis, 2023) vol. 202, pp. 8469–8488 ; <https://dl.acm.org/doi/10.5555/3618408.3618748>.

- 19* Octo Model Team, D. Ghosh, H. Walke, K. Pertsch, K. Black, O. Mees, S. Dasari, J. Hejna, T. Kreiman, C. Xu, J. Luo, Y. L. Tan, LY Chen, P. Sanketi, Q. Vuong, T. Xiao, D. Sadigh,... S. Levine, Octo : Une politique de robot généraliste open source, arXiv [cs.RO] (2024) ; <http://arxiv.org/abs/2405.12213>.
- 20 M. Firat, S. Kuleli , Et si GPT4 devenait autonome : le projet Auto-GPT et les cas d'utilisation. *Journal of Emerging Technologies informatiques* 3, 1–6 (2024) ; <https://doi.org/10.57020/ject.1297961>.
- 21* Y. Wang, T. Shen, L. Liu, J. Xie, Sibyl : Cadre d'agent simple mais efficace pour le raisonnement complexe du monde réel, arXiv [cs.AI] (2024) ; <http://arxiv.org/abs/2407.10718>.
- 22* C. Lu, C. Lu, RT Lange, J. Foerster, J. Clune, D. Ha, The AI Scientist : Vers une approche ouverte entièrement automatisée Découverte scientifique, arXiv [cs.AI] (2024) ; <http://arxiv.org/abs/2408.06292>.
- 23 J. Abramson, J. Adler, J. Dunger, R. Evans, T. Green, A. Pritzel, O. Ronneberger, L. Willmore, A. J. Ballard, J. Bambrick, SW Bodenstern, DA Evans, C.-C. Hung, M. O'Neill, D. Reiman, K. Tunyasuvunakool, Z. Wu,... JM Jumper, Prédiction précise de la structure des interactions biomoléculaires avec AlphaFold 3. *Nature* 630, 493–500 (2024) ; <https://doi.org/10.1038/s41586-024-07487-w>.
- 24 Y. LeCun, Y. Bengio, G. Hinton, Deep Learning. *Nature* 521, 436–444 (2015) ; <https://doi.org/10.1038/nature14539>.
- 25 A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, AN Gomez, Ł. U. Kaiser, I. Polosukhin, « L'attention est tout ce dont vous avez besoin » dans *Advances in Neural Information Processing Systems (NIPS 2017)* (Curran Associates, Inc., 2017) vol. 30 ; https://papers.nips.cc/paper_files/paper/2017/hash/3f5ee243547dee91fbd053c1c4a845aa-Abstract.html.
- 26 J. Sevilla, L. Heim, A. Ho, T. Besiroglu, M. Hobbhahn, P. Villalobos, « Compute Trends Across Three Eras of Machine Learning » dans *Conférence internationale conjointe sur les réseaux neuronaux 2022 (IJCNN 2022)* (Padoue, Italie, 2022), pp. 1–8 ; <https://doi.org/10.1109/IJCNN55064.2022.9891914>.
- 27 B. Cottier, R. Rahman, L. Fattorini, N. Maslej, D. Owen, Combien coûte la formation des modèles Frontier AI ?, *Epoch AI* (2024) ; <https://epochai.org/blog/how-much-does-it-cost-to-train-frontier-ai-models>.
- 28 C. Zhou, P. Liu, P. Xu, S. Iyer, J. Sun, Y. Mao, X. Ma, A. Efrat, P. Yu, L. Yu, S. Zhang, G. Ghosh, M. Lewis, L. Zettlemoyer, O. Levy, « LIMA : moins c'est plus pour l'alignement » dans *37e Conférence sur les systèmes de traitement de l'information neuronale (NeurIPS 2023)* (La Nouvelle-Orléans, LA, États-Unis, 2023) ; <https://openreview.net/forum?id=KBMOkMx2he>.
- 29 R. Rafailov, A. Sharma, E. Mitchell, CD Manning, S. Ermon, C. Finn, « Optimisation des préférences directes : votre modèle de langage est secrètement un modèle de récompense » dans *37e Conférence sur les systèmes de traitement de l'information neuronale (NeurIPS 2023)* (La Nouvelle-Orléans, LA, États-Unis, 2023) ; <https://openreview.net/forum?id=HPuSIXJaa9>.
- 30 L. Ouyang, J. Wu, X. Jiang, D. Almeida, C. Wainwright, P. Mishkin, C. Zhang, S. Agarwal, K. Slama, A. Gray, J. Schulman, J. Hilton, F. Kelton, L. Miller, M. Simens, A. Askell, P. Welinder, ... R. Lowe, « Entraîner des modèles linguistiques à suivre des instructions avec un retour humain » dans *36e Conférence sur les systèmes de traitement de l'information neuronale (NeurIPS 2022)* (La Nouvelle-Orléans, LA, États-Unis, 2022) ; <https://openreview.net/forum?id=TG8KACxEON>.
- 31* Y. Bai, A. Jones, K. Ndousse, A. Askell, A. Chen, N. DasSarma, D. Drain, S. Fort, D. Ganguli, T. Henighan, N. Joseph, S. Kadavath, J. Kernion, T. Conerly, S. El-Showk, N. Elhage, Z. Hatfield-Dodds, ... J. Kaplan, Former un assistant utile et inoffensif avec l'apprentissage par renforcement à partir de commentaires humains, arXiv [cs.CL] (2022) ; <http://arxiv.org/abs/2204.05862>.
- 32* N. McAleese, RM Pokorny, JFC Uribe, E. Nitishinskaya, M. Trebacz, J. Leike, LLM Critics Help Catch LLM Bugs, arXiv [cs.SE] (2024) ; <http://arxiv.org/abs/2407.00215>.
- 33* H. Lee, S. Phatale, H. Mansoor, T. Mesnard, J. Ferret, K. Lu, C. Bishop, E. Hall, V. Carbune, A. Rastogi, S. Prakash, RLAIF : Mise à l'échelle de l'apprentissage par renforcement à partir du feedback humain avec le feedback de l'IA, arXiv [cs.CL] (2023) ; <http://arxiv.org/abs/2309.00267>.
- 34 M. Mitchell, S. Wu, A. Zaldivar, P. Barnes, L. Vasserman, B. Hutchinson, E. Spitzer, ID Raji, T. Gebru, « Cartes modèles pour un reporting modèle » dans *Actes de la Conférence sur l'équité, la responsabilité et la transparence (FAT* '19)* (Association for Computing Machinery, New York, NY, États-Unis, 2019), pp. 220–229 ; <https://doi.org/10.1145/3287560.3287596>.
- 35* I. Solaiman, Le gradient de diffusion de l'IA générative : méthodes et considérations, arXiv [cs.CY] (2023) ; <http://arxiv.org/abs/2302.04844>.
- 36* Open Source Initiative, La définition de l'IA Open Source – 1.0-RC2, Open Source Initiative (2024) ; <https://opensource.org/ai/drafts/the-open-source-ai-definition-1-0-rc2>.
- 37* A. Dubey, A. Jauhri, A. Pandey, A. Kadian, A. Al-Dahle, A. Letman, A. Mathur, A. Schelten, A. Yang, A. Fan, A. Goyal, A. Hartshorn, A. Yang, A. Mitra, A. Srivankumar, A. Korenev, A. Hinsvark,... Z. Zhao, « Le troupeau de lamas 3 de Modèles » (Méta, 2024) ; <https://ai.meta.com/research/publications/the-llama-3-herd-of-models/>.
- 38 M. Stein, C. Dunlop, Safe beyond Sale : surveillance post-déploiement de l'IA (2024) ; <https://www.adalovelaceinstitute.org/blog/post-deployment-monitoring-of-ai/>.

- 39 E. Shayegani, MA Al Mamun, Y. Fu, P. Zaree, Y. Dong, N. Abu-Ghazaleh, Enquête sur les vulnérabilités dans les grands modèles linguistiques révélées par des attaques contradictoires, arXiv [cs.CL] (2023) ; <http://arxiv.org/abs/2310.10844>.
- 40 RT McCoy, S. Yao, D. Friedman, MD Hardy, TL Griffiths, Lorsqu'un modèle de langage est optimisé pour le raisonnement, présente-t-il encore des signes d'autorégression ? Une analyse d'OpenAI o1, arXiv [cs.CL] (2024) ; <http://arxiv.org/abs/2410.01792>.
- 41 U. Anwar, A. Saparov, J. Rando, D. Paleka, M. Turpin, P. Hase, ES Lubana, E. Jenner, S. Casper, O. Sourbut, BL Edelman, Z. Zhang, M. Günther, A. Korinek, J. Hernandez-Orallo, L. Hammond, E. Bigelow, ... D. Krueger, Défis fondamentaux pour assurer l'alignement et la sécurité des grands modèles de langage, arXiv [cs.LG] (2024) ; <http://arxiv.org/abs/2404.09932>.
- 42* RT McCoy, S. Yao, D. Friedman, M. Hardy, TL Griffiths, Embers of Autoregression : comprendre les grands modèles de langage à travers le problème qu'ils sont formés à résoudre, arXiv [cs.CL] (2023) ; <http://arxiv.org/abs/2309.13638>.
- 43 Y. Razeghi, RL Logan IV, M. Gardner, S. Singh, Impact des fréquences de termes de pré-entraînement sur le raisonnement à quelques coups, arXiv [cs.CL] (2022) ; <http://arxiv.org/abs/2202.07206>.
- 44* T. Shevlane, S. Farquhar, B. Garfinkel, M. Phuong, J. Whittlestone, J. Leung, D. Kokotajlo, N. Marchal, M. Anderjung, N. Kolt, L. Ho, D. Siddarth, S. Avin, W. Hawkins, B. Kim, I. Gabriel, V. Bolina, ... A. Dafoe, « Évaluation des modèles pour les risques extrêmes » (Google DeepMind, 2023) ; <http://arxiv.org/abs/2305.15324>.
- 45 R. Bommasani, D. Soylu, TI Liao, KA Creel, P. Liang, Ecosystem Graphs : The Social Footprint of Foundation Modèles, arXiv [cs.LG] (2023) ; <http://arxiv.org/abs/2303.15772>.
- 46* A. Das, W. Kong, R. Sen, Y. Zhou, Un modèle de base basé uniquement sur un décodeur pour la prévision des séries chronologiques, arXiv [cs.CL] (2023) ; <http://arxiv.org/abs/2310.10688>.
- 47* P. Dhariwal, H. Jun, C. Payne, JW Kim, A. Radford, I. Sutskever, « Jukebox : un modèle génératif pour la musique » (OpenAI, 2020) ; <http://arxiv.org/abs/2005.00341>.
- 48* H. Touvron, L. Martin, K. Stone, P. Albert, A. Almahairi, Y. Babaei, N. Bashlykov, S. Batra, P. Bhargava, S. Bhosale, D. Bikel, L. Blecher, CC Ferrer, M. Chen, G. Cucurull, D. Esiobu, J. Fernandes, ... T. Scialom, « Llama 2 : Fondation ouverte et modèles de chat affinés » (Meta AI, 2023) ; <http://arxiv.org/abs/2307.09288>.
- 49* Équipe Gemini, P. Georgiev, VI Lei, R. Burnell, L. Bai, A. Gulati, G. Tanzer, D. Vincent, Z. Pan, S. Wang, S. Mariooryad, Y. Ding, X. Geng, F. Alcober, R. Frostig, M. Omernick, L. Walker, ... O. Vinyals, « Gemini 1.5 : Débloquer la compréhension multimodale à travers des millions de jetons de contexte » (Google DeepMind, 2024) ; https://storage.googleapis.com/deepmind-media/gemini/gemini_v1_5_report.pdf.
- 50* Anthropic, « La famille modèle Claude 3 : Opus, Sonnet, Haiku » (Anthropic, 2024) ; https://www-cdn.anthropic.com/de8ba9b01c9ab7cbabf5c33b80b7bbc618857627/Model_Card_Claude_3.pdf.
- 51* OpenAI, « Carte système GPT-4 » (OpenAI, 2023) ; <https://cdn.openai.com/papers/gpt-4-system-card.pdf>.
- 52* AQ Jiang, A. Sablayrolles, A. Roux, A. Mensch, B. Savary, C. Bamford, DS Chaplot, D. de las Casas, EB Hanna, F. Bressand, G. Lengyel, G. Bour, G. Lampe, LR Lavaud, L. Saulnier, M.-A. Lachaux, P. Stock, ... WE Sayed, Mixtral d'experts, arXiv [cs.LG] (2024) ; <http://arxiv.org/abs/2401.04088>.
- 53* A. Yang, B. Yang, B. Hui, B. Zheng, B. Yu, C. Zhou, C. Li, C. Li, D. Liu, F. Huang, G. Dong, H. Wei, H. Lin, J. Tang, J. Wang, J. Yang, J. Tu, ... Z. Fan, rapport technique Qwen2, arXiv [cs.CL] (2024) ; <http://arxiv.org/abs/2407.10671>.
- 54* DeepSeek-AI, A. Liu, B. Feng, B. Wang, B. Wang, B. Liu, C. Zhao, C. Dengr, C. Ruan, D. Dai, D. Guo, D. Yang, D. Chen, D. Ji, E. Li, F. Lin, F. Luo, ... Z. Xie, DeepSeek-V2 : un modèle de langage à base de mélange d'experts puissant, économique et efficace, arXiv [cs.CL] (2024) ; <http://arxiv.org/abs/2405.04434>.
- 55 L. Wang, C. Ma, X. Feng, Z. Zhang, H. Yang, J. Zhang, Z. Chen, J. Tang, X. Chen, Y. Lin, WX Zhao, Z. Wei, J. Wen, A. Enquête sur les agents autonomes basés sur des modèles de langage de grande taille. *Frontiers of Computer Science* 18, 186345 (2024) ; <https://doi.org/10.1007/s11704-024-40231-1>.
- 56 A. Fan, B. Gokkaya, M. Harman, M. Lyubarskiy, S. Sengupta, S. Yoo, JM Zhang, « Grands modèles de langage pour « Ingénierie logicielle : enquête et problèmes ouverts » dans la Conférence internationale IEEE/ACM 2023 sur l'ingénierie logicielle : L'avenir de l'ingénierie logicielle (ICSE-FoSE) (2023), pp. 31–53 ; <https://doi.org/10.1109/ICSE-FoSE59343.2023.00008>.
- 57* S. Chen, S. Liu, L. Zhou, Y. Liu, X. Tan, J. Li, S. Zhao, Y. Qian, F. Wei, VALL-E 2 : Les modèles de langage de codec neuronal sont des synthétiseurs de texte à parole à parité humaine Zero-Shot, arXiv [cs.CL] (2024) ; <http://arxiv.org/abs/2406.05370>.
- 58* OpenAI, « Carte système GPT-4V(ision) » (OpenAI, 2023) ; <https://cdn.openai.com/gpt-4o-system-card.pdf>.
- 59* P. Agrawal, S. Antoniak, EB Hanna, B. Bout, D. Chaplot, J. Chudnovsky, D. Costa, B. De Monicault, S. Garg, T. Gervet, S. Ghosh, A. Héliou, P. Jacob, AQ Jiang, K. Khandelwal, T. Lacroix, G. Lample, ... S. Yang, Pixtral 12B, arXiv [cs.CV] (2024) ; <http://arxiv.org/abs/2410.07073>.

- 60* P. Wang, S. Bai, S. Tan, S. Wang, Z. Fan, J. Bai, K. Chen, X. Liu, J. Wang, W. Ge, Y. Fan, K. Dang, M. Du, X. Ren, R. Men, D. Liu, C. Zhou, ... J. Lin, Qwen2-VL : Améliorer la perception du monde par le modèle vision-langage à n'importe quelle résolution, arXiv [cs.CV] (2024) ; <http://arxiv.org/abs/2409.12191>.
- 61 A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, N. Houlsby, « Une image vaut 16x16 mots : transformateurs pour la reconnaissance d'images à grande échelle » dans la 9e Conférence internationale sur les représentations d'apprentissage (ICLR 2021) (virtuelle, 2020) ; <https://openreview.net/forum?id=YicbFdNTTy>.
- 62* A. Kirillov, E. Mintun, N. Ravi, H. Mao, C. Rolland, L. Gustafson, T. Xiao, S. Whitehead, A.C. Berg, W.-Y. Lo, P. Dollár, R. Girshick, « Segmentez n'importe quoi » (Meta AI, 2023) ; <http://arxiv.org/abs/2304.02643>.
- 63* A. Bardes, Q. Garrido, J. Ponce, X. Chen, M. Rabbat, Y. LeCun, M. Assran, N. Ballas, « Revisiting Feature Prediction pour l'apprentissage des représentations visuelles à partir de vidéos » (Meta, 2024).
- 64* L'équipe Movie Gen, « Movie Gen : un ensemble de modèles de fondations médiatiques » (Meta, 2024) ; <https://ai.meta.com/static-resource/movie-gen-research-paper>.
- 65* J. Liang, W. Huang, F. Xia, P. Xu, K. Hausman, B. Ichter, P. Florence, A. Zeng, « Le code comme politique : modèle de langage « Programmes pour le contrôle incarné » dans l'atelier sur le langage et la robotique au CoRL 2022 (2022) ; <https://openreview.net/forum?id=fmvtvpopfLC6>.
- 66 B. Ichter, A. Brohan, Y. Chebotar, C. Finn, K. Hausman, A. Herzog, D. Ho, J. Ibarz, A. Irpan, E. Jang, R. Julian, D. Kalashnikov, S. Levine, Y. Lu, C. Parada, K. Rao, P. Sermanet, ... CK Fu, « Fais ce que je peux, pas ce que je dis : ancrer le langage dans les capacités robotiques » dans les actes de la 6e conférence annuelle sur l'apprentissage des robots (CoRL) (PMLR, Auckland, Nouvelle-Zélande, 2022) vol. 205 ; https://openreview.net/forum?id=bdHkMjBJG_w.
- 67 Collaboration ouverte X-Embodiment, A. O'Neill, A. Rehman, A. Maddukuri, A. Gupta, A. Padalkar, A. Lee, A. Pooley, A. Gupta, A. Mandlekar, A. Jain, A. Tung, A. Bewley, A. Herzog, A. Irpan, A. Khazatsky, A. Rai, ... Z. Lin, Open X-Embodiment : ensembles de données d'apprentissage robotique et modèles RT-X, arXiv [cs.RO] (2023) ; <http://arxiv.org/abs/2310.08864>.
- 68* J.-J. Hwang, R. Xu, H. Lin, W.-C. Hung, J. Ji, K. Choi, D. Huang, T. He, P. Covington, B. Sapp, Y. Zhou, J. Guo, D. Anguelov, M. Tan, EMMA : Modèle multimodal de bout en bout pour la conduite autonome, arXiv [cs.CV] (2024) ; <http://arxiv.org/abs/2410.23262>.
- 69 R. Firoozi, J. Tucker, S. Tian, A. Majumdar, J. Sun, W. Liu, Y. Zhu, S. Song, A. Kapoor, K. Hausman, B. Ichter, D. Driess, J. Wu, C. Lu, M. Schwager, Modèles de base en robotique : applications, défis et avenir, arXiv [cs.RO] (2023) ; <http://arxiv.org/abs/2312.07843>.
- 70 H. Fang, H. Fang, Z. Tang, J. Liu, J. Wang, H. Zhu, C. Lu, RH20T : Un ensemble complet de données robotiques pour l'apprentissage de diverses compétences en une seule fois. Conférence internationale de l'IEEE sur la robotique et l'automatisation, 653–660 (2023) ; <https://doi.org/10.1109/ICRA57147.2024.10611615>.
- 71 A. Khazatsky, K. Pertsch, S. Nair, A. Balakrishna, S. Dasari, S. Karamcheti, S. Nasiriany, MK Srirama, LY Chen, K. Ellis, PD Fagan, J. Hejna, M. Itkina, M. Lepert, YJ Ma, PT Miller, J. Wu, ... C. Finn, DROID : un ensemble de données de manipulation de robots à grande échelle dans la nature, arXiv [cs.RO] (2024) ; <http://arxiv.org/abs/2403.12945>.
- 72 J. Wang, Z. Wu, Y. Li, H. Jiang, P. Shu, E. Shi, H. Hu, C. Ma, Y. Liu, X. Wang, Y. Yao, X. Liu, H. Zhao, Z. Liu, H. Dai, L. Zhao, B. Ge, ... S. Zhang, Grands modèles de langage pour la robotique : opportunités, défis et perspectives, arXiv [cs.RO] (2024) ; <http://arxiv.org/abs/2401.04334>.
- 73* Chai Discovery, J. Boiteaud, J. Dent, M. McPartlon, J. Meier, V. Reis, A. Rogozhnikov, K. Wu, Chai-1 : Décoder les Interactions moléculaires de la vie, bioRxiv [préimpression] (2024) ; <https://doi.org/10.1101/2024.10.10.615955>.
- 74 R. Bommasani, DA Hudson, E. Adeli, R. Altman, S. Arora, S. von Arx, MS Bernstein, J. Bohg, A. Bosselut, E. Brunskill, E. Brynjolfsson, S. Buch, D. Card, R. Castellon, N. Chatterji, A. Chen, K. Creel, ... P. Liang, Sur les opportunités et les risques des modèles de fondation, arXiv [cs.LG] (2021) ; <http://arxiv.org/abs/2108.07258>.
- 75 P. Bryant, G. Pozzati, A. Elofsson, Prédiction améliorée des interactions protéine-protéine à l'aide d'AlphaFold2. Nature Communications 13, 1265 (2022) ; <https://doi.org/10.1038/s41467-022-28865-w>.
- 76 A. Madani, B. Krause, ER Greene, S. Subramanian, BP Mohr, JM Holton, JL Olmos, C. Xiong, ZZ Sun, R. Socher, JS Fraser, N. Naik, Les grands modèles de langage génèrent des séquences de protéines fonctionnelles dans diverses familles. Nature Biotechnology 41, 1099–1106 (2023) ; <https://doi.org/10.1038/s41587-022-01618-2>.
- 77 T. Davidson, J.-S. Denain, P. Villalobos, G. Bas, « Les capacités de l'IA peuvent être considérablement améliorées sans recyclage coûteux » (Epoch AI, 2023) ; <http://arxiv.org/abs/2312.07413>.
- 78 G. Mialon, R. Dessi, M. Lomeli, C. Nalmpantis, R. Pasunuru, R. Raileanu, B. Roziere, T. Schick, J. Dwivedi-Yu, A. Celikyilmaz, E. Grave, Y. LeCun, T. Scialom, Modèles de langage augmentés : une enquête. Transactions on Machine Learning Research (2023) ; <https://openreview.net/pdf?id=jh7wH2AzKK>.

- 79* X. Wang, J. Wei, D. Schuurmans, Q. Le, E. Chi, S. Narang, A. Chowdhery, D. Zhou, L'auto-cohérence améliore la chaîne de raisonnement dans les modèles de langage, arXiv [cs.CL] (2022) ; <http://arxiv.org/abs/2203.11171>.
- 80* B. Brown, J. Juravsky, R. Ehrlich, R. Clark, QV Le, C. Ré, A. Mirhoseini, Grands singes du langage : mise à l'échelle du calcul d'inférence avec échantillonnage répété, arXiv [cs.LG] (2024) ; <http://arxiv.org/abs/2407.21787>.
- 81 S. Yao, D. Yu, J. Zhao, I. Shafran, TL Griffiths, Y. Cao, KR Narasimhan, « Arbre de pensées : résolution délibérée de problèmes avec de grands modèles de langage » dans 37e Conférence sur les systèmes de traitement de l'information neuronale (NeurIPS 2023) (La Nouvelle-Orléans, LA, États-Unis, 2023) ; <https://openreview.net/forum?id=5Xc1ecxO1h>.
- 82 T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger, T. Henighan, R. Child, A. Ramesh, D. Ziegler, ... D. Amodei, « Les modèles linguistiques sont des apprenants peu nombreux » dans Advances in Neural Information Processing Systems (Curran Associates, Inc., 2020) vol. 33, pp. 1877–1901 ; <https://papers.nips.cc/paper/2020/hash/1457c0d6bfc4967418bfb8ac142f64a-Abstract.html>.
- 83 J. Wei, X. Wang, D. Schuurmans, M. Bosma, B. Ichter, F. Xia, E. Chi, QV Le, D. Zhou, « L'incitation à la chaîne de pensée suscite le raisonnement dans les grands modèles linguistiques » dans Advances in Neural Information Processing Systems (NeurIPS 2022) (La Nouvelle-Orléans, LA, États-Unis, 2022) vol. 35, pp. 24824–24837 ; https://proceedings.neurips.cc/paper_files/paper/2022/hash/9d5609613524ecf4f15af0f7b31abca4-Abstract-Conference.html.
- 84 T. Kojima, SS Gu, M. Reid, Y. Matsuo, Y. Iwasawa, « Les grands modèles de langage sont des raisonneurs à coup zéro » dans NeurIPS (La Nouvelle-Orléans, LA, États-Unis, 2022) ; http://papers.nips.cc/paper_files/paper/2022/hash/8bb0d291acd4ac06ef112099c16f326-Abstract-Conference.html.
- 85* R. Nakano, J. Hilton, S. Balaji, J. Wu, L. Ouyang, C. Kim, C. Hesse, S. Jain, V. Kosaraju, W. Saunders, X. Jiang, K. Cobbe, T. Eloundou, G. Krueger, K. Button, M. Knight, B. Chess, J. Schulman, « WebGPT : réponse aux questions assistée par navigateur avec retour humain » (OpenAI, 2021) ; <http://arxiv.org/abs/2112.09332>.
- 86* L. Gao, A. Madaan, S. Zhou, U. Alon, P. Liu, Y. Yang, J. Callan, G. Neubig, PAL : Modèles de langage assistés par programme, arXiv [cs.CL] (2022) ; <https://doi.org/10.48550/arXiv.2211.10435>.
- 87 I. Drori, S. Zhang, R. Shuttleworth, L. Tang, A. Lu, E. Ke, K. Liu, L. Chen, S. Tran, N. Cheng, R. Wang, N. Singh, TL Patti, J. Lynch, A. Shporer, N. Verma, E. Wu, G. Strang, Un réseau neuronal résout, explique et génère des problèmes mathématiques universitaires par synthèse de programme et apprentissage en quelques coups au niveau humain, arXiv [cs.LG] (2021) ; <https://pnas.org/doi/full/10.1073/pnas.2123433119>.
- 88* W. Chen, X. Ma, X. Wang, WW Cohen, Programme d'incitation à la pensée : démêler le calcul de Raisonnement pour les tâches de raisonnement numérique, arXiv [cs.CL] (2022) ; <http://arxiv.org/abs/2211.12588>.
- 89 W. Huang, P. Abbeel, D. Pathak, I. Mordatch, Modèles de langage comme planificateurs à coup zéro : extraction de connaissances exploitables pour les agents incarnés. (2022) ; <https://openreview.net/forum?id=6NT1a56mNim>.
- 90 I. Dasgupta, C. Kaeser-Chen, K. Marino, A. Ahuja, S. Babayan, F. Hill, R. Fergus, « Collaboration avec des modèles linguistiques pour le raisonnement incarné » dans Deuxième atelier sur l'apprentissage du langage et du renforcement (2022) ; <https://openreview.net/forum?id=YoS-abmWJc>.
- 91 Epoch AI, tableau de bord d'analyse comparative de l'IA (2024) ; <https://epoch.ai/data/ai-benchmarking-dashboard>.
- 92* OpenAI, Apprendre à raisonner avec les LLM (2024) ; <https://openai.com/index/learning-to-reason-with-llms/>.
- 93 P. Villalobos, D. Atkinson, « Échange de calculs entre formation et inférence » (Epoch AI, 2023) ; <https://epochai.org/blog/trading-off-compute-in-training-and-inference>.
- 94* C. Snell, J. Lee, K. Xu, A. Kumar, La mise à l'échelle optimale du temps de calcul du test LLM peut être plus efficace que la mise à l'échelle Paramètres du modèle, arXiv [cs.LG] (2024) ; <http://arxiv.org/abs/2408.03314>.
- 95 X. Hu, J. Chen, X. Li, Y. Guo, L. Wen, PS Yu, Z. Guo, Les grands modèles de langage connaissent-ils les faits ?, arXiv [cs.CL] (2023) ; <http://arxiv.org/abs/2310.05177>.
- 96 R. Xu, Z. Qi, Z. Guo, C. Wang, H. Wang, Y. Zhang, W. Xu, « Conflits de connaissances pour les LLM : une enquête » dans Actes de la Conférence 2024 sur les méthodes empiriques en traitement du langage naturel (Association for Computational Linguistics, Stroudsburg, PA, États-Unis, 2024), pp. 8541–8565 ; <https://doi.org/10.18653/v1/2024.emnlp-main.486>.
- 97 M. Turpin, J. Michael, E. Perez, SR Bowman, « Les modèles linguistiques ne disent pas toujours ce qu'ils pensent : explications infidèles dans l'incitation à la chaîne de pensée » dans 37e Conférence sur les systèmes de traitement de l'information neuronale (NeurIPS 2023) (La Nouvelle-Orléans, LA, États-Unis, 2023) ; <https://openreview.net/forum?id=bzs4uPLXvi>.
- 98 M. Sharma, M. Tong, T. Korbak, D. Duvenaud, A. Askell, SR Bowman, E. Durmus, Z. Hatfield-Dodds, SR Johnston, SM Kravec, T. Maxwell, S. McCandlish, K. Ndousse, O. Rausch, N. Schiefer, D. Yan, M. Zhang, E. Perez, « Vers une compréhension de la flagornerie dans les modèles linguistiques » dans la 12e Conférence internationale sur l'apprentissage

Représentations (ICLR 2024) (Vienne, Autriche, 2023) ; <https://openreview.net/forum?id=tvhaxkMKAn>.

- 99* Z. Wu, L. Qiu, A. Ross, E. Akyürek, B. Chen, B. Wang, N. Kim, J. Andreas, Y. Kim, Raisonnement ou récitation ? Exploration de la Capacité et limites des modèles linguistiques à travers des tâches contrefactuelles, arXiv [cs.CL] (2023) ; <http://arxiv.org/abs/2307.02477>.
- 100 L. Zhang, X. Zhai, Z. Zhao, Y. Zong, X. Wen, B. Zhao, « Et si la télévision était éteinte ? Examen des capacités de raisonnement contrefactuel des modèles de langage multimodaux » dans Actes de la conférence IEEE/CVF sur la vision par ordinateur et la reconnaissance de formes (2024), pp. 21853–21862 ; https://openaccess.thecvf.com/content/CVPR2024/papers/Zhang_What_If_the_TV_Was_Off_Examining_Cou_Capacités_de_raisonnement_interfactuel_CVPR_2024_paper.pdf.
- 101 Z. Ji, N. Lee, R. Frieske, T. Yu, D. Su, Y. Xu, E. Ishii, YJ Bang, A. Madotto, P. Fung, Survey of Hallucination in Natural Génération de langage. ACM Computing Surveys 55, 1–38 (2023) ; <https://doi.org/10.1145/3571730>.
- 102* Y. Zhang, Y. Li, L. Cui, D. Cai, L. Liu, T. Fu, X. Huang, E. Zhao, Y. Zhang, Y. Chen, L. Wang, AT Luu, W. Bi, F. Shi, S. Shi, Le chant de la sirène dans l'océan de l'IA : une enquête sur les hallucinations dans les grands modèles linguistiques, arXiv [cs.CL] (2023) ; <http://arxiv.org/abs/2309.01219>.
- 103 M. Zhang, O. Press, W. Merrill, A. Liu, NA Smith, Comment les hallucinations du modèle de langage peuvent faire boule de neige, arXiv [cs.CL] (2023) ; <http://arxiv.org/abs/2305.13534>.
- 104 L. Huang, W. Yu, W. Ma, W. Zhong, Z. Feng, H. Wang, Q. Chen, W. Peng, X. Feng, B. Qin, T. Liu, Une enquête sur Hallucination dans les grands modèles linguistiques : principes, taxonomie, défis et questions ouvertes, arXiv [cs.CL] (2023) ; <http://arxiv.org/abs/2311.05232>.
- 105 V. Rawte, A. Sheth, A. Das, Une étude des hallucinations dans les modèles de grandes fondations, arXiv [cs.AI] (2023) ; <http://arxiv.org/abs/2309.05922>.
- 106 J. Liu, W. Wang, D. Wang, N. Smith, Y. Choi, H. Hajishirzi, « Vera : un modèle d'estimation de plausibilité à usage général pour les énoncés de sens commun » dans Actes de la Conférence 2023 sur les méthodes empiriques en traitement du langage naturel, H. Bouamor, J. Pino, K. Bali, éd. (Association for Computational Linguistics, Singapour, 2023), pp. 1264–1287 ; <https://doi.org/10.18653/v1/2023.emnlp-main.81>.
- 107 A. Leidingger, R. Van Rooij, E. Shutova, « Les LLM sont-ils des raisonneurs classiques ou non monotones ? Leçons tirées des génériques » dans Actes de la 62e réunion annuelle de l'Association for Computational Linguistics (Volume 2 : articles courts), L.-W. Ku, A. Martins, V. Srikumar, éd. (Association for Computational Linguistics, Bangkok, Thaïlande, 2024) ; <https://doi.org/10.18653/v1/2024.acl-short.51>.
- 108 M. Mitchell, Le défi de l'IA pour comprendre le monde. Science 382, eadm8175 (2023) ; <https://doi.org/10.1126/science.adm8175>.
- 109 D. Halawi, F. Zhang, C. Yueh-Han, J. Steinhardt, Approche de la prévision au niveau humain avec des modèles linguistiques, arXiv [cs.LG] (2024) ; <http://arxiv.org/abs/2402.18563>.
- 110* I. Mirzadeh, K. Alizadeh, H. Shahrokhi, O. Tuzel, S. Bengio, M. Farajtabar, GSM-Symbolic : comprendre les limites du raisonnement mathématique dans les grands modèles de langage, arXiv [cs.LG] (2024) ; <http://arxiv.org/abs/2410.05229>.
- 111* F. Shi, X. Chen, K. Misra, N. Scales, D. Dohan, EH Chi, N. Schärli, D. Zhou, « Les grands modèles linguistiques peuvent être facilement distraits par un contexte non pertinent » dans Actes de la 40e Conférence internationale sur l'apprentissage automatique (PMLR, 2023), pp. 31210–31227 ; <https://proceedings.mlr.press/v202/shi23a.html>.
- 112* A. Hosseini, A. Sordani, D. Toyama, A. Courville, R. Agarwal, Tous les raisonneurs LLM ne sont pas égaux, arXiv [cs.LG] (2024) ; <http://arxiv.org/abs/2410.01748>.
- 113 KZ Cui, M. Demirer, S. Jaffe, L. Musolff, S. Peng, T. Salz, Les effets de l'IA générative sur la productivité : résultats d'une expérience sur le terrain avec GitHub Copilot. Une exploration de l'IA générative par le MIT (2024) ; <https://mit-genai.pubpub.org/pub/v5iixksv/release/2>.
- 114* S. Peng, E. Kalliamvakou, P. Cihon, M. Demirer, L'impact de l'IA sur la productivité des développeurs : données de GitHub Copilote, arXiv [cs.SE] (2023) ; <https://www.semanticscholar.org/reader/038f249ab708cebae2a58265b768b9b1cbadad3a>.
- 115 A. Ziegler, E. Kalliamvakou, XA Li, A. Rice, D. Rifkin, S. Simister, G. Sittampalam, E. Aftandilian, Mesure de l'impact de GitHub Copilot sur la productivité. Communications de l'ACM 67, 54–63 (2024) ; <https://doi.org/10.1145/3633453>.
- Enquête 2024 auprès des développeurs Stack Overflow (2024) ; <https://survey.stackoverflow.co/2024/>.
- 117 Enquête auprès des développeurs Stack Overflow 2023, Stack Overflow (2023) ; https://survey.stackoverflow.co/2023/?utm_source=social-share&utm_medium=social&utm_campaign=dev-survey-2023.

- 118 X. Liu, H. Yu, H. Zhang, Y. Xu, X. Lei, H. Lai, Y. Gu, H. Ding, K. Men, K. Yang, S. Zhang, X. Deng, A. Zeng, Z. Du, C. Zhang, S. Shen, T. Zhang, ... J. Tang, AgentBench : évaluation des LLM en tant qu'agents, arXiv [cs.AI] (2023) ; <http://arxiv.org/abs/2308.03688>.
- 119 S. Yao, H. Chen, J. Yang, K. Narasimhan, WebShop : vers une interaction Web évolutive dans le monde réel avec des agents linguistiques ancrés, arXiv [cs.CL] (2022) ; <http://arxiv.org/abs/2207.01206>.
- 120 AM Bran, S. Cox, O. Schilter, C. Baldassari, A. White, P. Schwaller, « Augmentation des grands modèles de langage avec « Outils de chimie » dans la 37e conférence sur les systèmes de traitement de l'information neuronale (NeurIPS 2023) Atelier sur l'IA pour la science (La Nouvelle-Orléans, LA, États-Unis, 2023) ; <https://openreview.net/forum?id=wdGIL6ix3l>.
- 121* AM Bran, S. Cox, O. Schilter, C. Baldassari, AD White, P. Schwaller, ChemCrow : Augmentation du langage à grande échelle Modèles avec des outils de chimie, arXiv [physics.chem-ph] (2023) ; <http://arxiv.org/abs/2304.05376>.
- 122 CE Jimenez, J. Yang, A. Wettig, S. Yao, K. Pei, O. Press, KR Narasimhan, « SWE-Bench : les modèles de langage peuvent-ils résoudre les problèmes Github du monde réel ? » dans 12e Conférence internationale sur les représentations d'apprentissage (2023) ; <https://openreview.net/pdf?id=VTF8yNQM66>.
- 123 *L. Jing, Z. Huang, X. Wang, W. Yao, W. Yu, K. Ma, H. Zhang, X. Du, D. Yu, DSBench : Dans quelle mesure les agents de science des données sont-ils en passe de devenir des experts en science des données ?, arXiv [cs.AI] (2024) ; <http://arxiv.org/abs/2409.07703>.
- 124 Z. Chen, S. Chen, Y. Ning, Q. Zhang, B. Wang, B. Yu, Y. Li, Z. Liao, C. Wei, Z. Lu, V. Dey, M. Xue, FN Baker, B. Burns, D. Adu-Ampratwum, X. Huang, X. Ning, ... H. Sun, ScienceAgentBench : Vers une évaluation rigoureuse des agents linguistiques pour la découverte scientifique basée sur les données, arXiv [cs.CL] (2024) ; <http://arxiv.org/abs/2410.05080>.
- 125* JS Chan, N. Chowdhury, O. Jaffe, J. Aung, D. Sherburn, E. Mays, G. Starace, K. Liu, L. Maksin, T. Patwardhan, L. Weng, A. Mądry, MLE-Bench : Évaluation des agents d'apprentissage automatique sur l'ingénierie de l'apprentissage automatique, arXiv [cs.CL] (2024) ; <http://arxiv.org/abs/2410.07095>.
- 126 Q. Huang, J. Vora, P. Liang, J. Leskovec, « MAgentBench : évaluation des agents linguistiques sur l'expérimentation de l'apprentissage automatique » dans Quarante et unième conférence internationale sur l'apprentissage automatique (2024) ; <https://openreview.net/pdf?id=1Fs1LvYQW>.
- 127 R. Fang, R. Bindu, A. Gupta, Q. Zhan, D. Kang, LLM Les agents peuvent pirater des sites Web de manière autonome, arXiv [cs.CR] (2024) ; <http://arxiv.org/abs/2402.06664>.
- 128 X. Liang, L. Ma, S. Guo, J. Han, H. Xu, S. Ma, X. Liang, CorNav : Agent autonome avec planification auto-corrective pour Navigation visuelle et linguistique Zero-Shot, arXiv [cs.CV] (2023) ; <http://arxiv.org/abs/2306.10322>.
- 129 METR, Détails sur l'évaluation préliminaire d'OpenAI o1-Preview par METR. (2024) ; <https://metr.github.io/autonomy-evals-guide/openai-o1-preview-report/>.
- 130* J. Yang, CE Jimenez, A. Wettig, K. Lieret, S. Yao, K. Narasimhan, O. Press, SWE-Agent : Agent-Ordinateur Les interfaces permettent l'ingénierie logicielle automatisée, arXiv [cs.SE] (2024) ; <http://arxiv.org/abs/2405.15793>.
- 131* CS Xia, Y. Deng, S. Dunn, L. Zhang, Agentless : démystifier les agents d'ingénierie logicielle basés sur LLM, arXiv [cs.SE] (2024) ; <http://arxiv.org/abs/2407.01489>.
- 132* X. Wang, B. Li, Y. Song, FF Xu, X. Tang, M. Zhuge, J. Pan, Y. Song, B. Li, J. Singh, HH Tran, F. Li, R. Ma, M. Zheng, B. Qian, Y. Shao, N. Muennighoff, ... G. Neubig, OpenHands : une plateforme ouverte pour les développeurs de logiciels d'IA en tant qu'agents généralistes, arXiv [cs.SE] (2024) ; <http://arxiv.org/abs/2407.16741>.
- 133* C.-L. Cheang, G. Chen, Y. Jing, T. Kong, H. Li, Y. Li, Y. Liu, H. Wu, J. Xu, Y. Yang, H. Zhang, M. Zhu, GR-2 : Un modèle génératif de vidéo-langage-action avec des connaissances à l'échelle du Web pour la manipulation de robots, arXiv [cs.RO] (2024) ; <http://arxiv.org/abs/2410.06158>.
- 134 B. Wang, J. Zhang, S. Dong, I. Fang, C. Feng, VLM See, Robot Do : vidéo de démonstration humaine vers plan d'action du robot via le modèle de langage visuel, arXiv [cs.RO] (2024) ; <http://arxiv.org/abs/2410.08792>.
- 135* S. Ye, J. Jang, B. Jeon, S. Joo, J. Yang, B. Peng, A. Mandlekar, R. Tan, Y.-W. Chao, BY Lin, L. Liden, K. Lee, J. Gao, L. Zettlemoyer, D. Fox, M. Seo, Pré-entraînement à l'action latente à partir de vidéos, arXiv [cs.RO] (2024) ; <http://arxiv.org/abs/2410.11758>.
- 136 M. Herrmann, FJD Lange, K. Eggenberger, G. Casalicchio, M. Wever, M. Feurer, D. Rügamer, E. Hüllermeier, A.-L. Boulesteix, B. Bischl, « Position : pourquoi nous devons repenser la recherche empirique en apprentissage automatique » dans Actes de la 41e Conférence internationale sur l'apprentissage automatique, R. Salakhutdinov, Z. Kolter, K. Heller, A. Weller, N. Oliver, J. Scarlett, F. Berkenkamp, éd. (PMLR, 2024) vol. 235, pp. 18228–18247 ; <https://proceedings.mlr.press/v235/herrmann24b.html>.
- 137 D. Hendrycks, C. Burns, S. Kadavath, A. Arora, S. Basart, E. Tang, D. Song, J. Steinhardt, « Mesure de la résolution de problèmes mathématiques avec l'ensemble de données MATH » dans 35e conférence sur les systèmes de traitement de l'information neuronale (NeurIPS 2021) Datasets and Benchmarks Track (Round 2) (Virtual, 2021) ; <https://openreview.net/forum?id=7Bywt2mQsCe>.

- 138 J. Au Yeung, Z. Kraljevic, A. Luintel, A. Balston, E. Idowu, RJ Dobson, JT Teo, Les chatbots IA ne sont pas encore prêts pour une utilisation clinique. *Frontiers in Digital Health* 5, 1161098 (2023) ; <https://doi.org/10.3389/fdgh.2023.1161098>.
- 139 D. Kiela, M. Bartolo, Y. Nie, D. Kaushik, A. Geiger, Z. Wu, B. Vidgen, G. Prasad, A. Singh, P. Ringshia, Z. Ma, T. Thrush, S. Riedel, Z. Waseem, P. Stenetorp, R. Jia, M. Bansal, ... A. Williams, « Dynabench : repenser l'analyse comparative en PNL » dans Actes de la conférence 2021 du chapitre nord-américain de l'Association for Computational Linguistics : Human Language Technologies (Association for Computational Linguistics, 2021), pp. 4110–4124 ; <https://doi.org/10.18653/v1/2021.naacl-main.324>.
- 140 D. Hendrycks, C. Burns, S. Basart, A. Zou, M. Mazeika, D. Song, J. Steinhardt, « Mesurer le multitâche massif « Compréhension du langage » dans la 9e Conférence internationale sur les représentations de l'apprentissage (ICLR 2021) (virtuelle, 2021) ; <https://openreview.net/forum?id=d7KBJmI3GmQ>.
- 141 D. Rein, BL Hou, AC Stickland, J. Petty, RY Pang, J. Dirani, J. Michael, SR Bowman, GPQA : une référence de questions et réponses à l'épreuve de Google au niveau supérieur, *arXiv [cs.AI]* (2023) ; <http://arxiv.org/abs/2311.12022>.
- 142 A. Srivastava, A. Rastogi, A. Rao, AAM Shob, A. Abid, A. Fisch, AR Brown, A. Santoro, A. Gupta, A. Garriga-Alonso, A. Kluska, A. Lewkowycz, A. Agarwal, A. Power, A. Ray, A. Warstadt, AW Kocurek, ... Z. Wu, Au-delà du jeu d'imitation : quantification et extrapolation des capacités des modèles de langage. *Transactions on Machine Learning Research* (2023) ; <https://openreview.net/forum?id=uyTL5Bvosj>.
- 143* L. Kilpatrick, SB Mallick, modèles Gemini prêts pour la production mis à jour, prix 1.5 Pro réduit, tarif augmenté Limites et plus encore, GEMINI (2024) ; <https://developers.googleblog.com/en/updated-gemini-models-reduced-15-pro-tarification-augmentation-des-limites-de-taux-et-plus/>.
- 144 M. Hobbhahn, L. Heim, G. Aydos, « Tendances en matière de matériel d'apprentissage automatique » (Epoch AI, 2023) ; <https://epochai.org/blog/trends-in-machine-learning-hardware>.
- 145* HW Chung, L. Hou, S. Longpre, B. Zoph, Y. Tay, W. Fedus, Y. Li, X. Wang, M. Dehghani, S. Brahma, A. Webson, SS Gu, Z. Dai, M. Suzgun, X. Chen, A. Chowdhery, A. Castro-Ros, ... J. Wei, Mise à l'échelle des modèles de langage affinés pour les instructions, *arXiv [cs.LG]* (2022) ; <http://arxiv.org/abs/2210.11416>.
- 146* OpenAI, GPT-4o Mini : faire progresser l'intelligence rentable (2024) ; <https://openai.com/index/gpt-4o-mini-faire-progresser-l'intelligence-rentable/>.
- 147* OpenAI, J. Achiam, S. Adler, S. Agarwal, L. Ahmad, I. Akkaya, FL Aleman, D. Almeida, J. Altschmidt, S. Altman, S. Anadkat, R. Avila, I. Babuschkin, S. Balaji, V. Balcom, P. Baltescu, H. Bao, ... B. Zoph, « Rapport technique GPT-4 » (OpenAI, 2024) ; <http://arxiv.org/abs/2303.08774>.
- 148* OpenAI, Tarification (2024) ; <https://openai.com/chatgpt/pricing/>.
- 149* Tarification ensemble, together.ai (2023) ; <https://www.together.ai/pricing>.
- 150 PAR Lin, Y. Deng, K. Chandu, F. Brahman, A. Ravichander, V. Pyatkin, N. Dziri, RL Bras, Y. Choi, WildBench : Analyse comparative des LLM avec des tâches difficiles auprès d'utilisateurs réels dans la nature, *arXiv [cs.CL]* (2024) ; <http://arxiv.org/abs/2406.04770>.
- 151* J. Wang, J. Wang, B. Athiwaratkun, C. Zhang, J. Zou, Le mélange d'agents améliore le modèle de langage à grande échelle Capacités, *arXiv [cs.CL]* (2024) ; <http://arxiv.org/abs/2406.04692>.
- 152 J. Sevilla, « La capacité de calcul d'entraînement des modèles d'IA de pointe augmente de 4 à 5 fois par an » (2024) ; <https://epoch.ai/blog/training-compute-of-frontier-ai-models-grows-by-4-5x-per-year>.
- 153 M. Mitchell, AB Palmarini, AK Moskvichev, « Comparaison des humains, GPT-4 et GPT-4V sur l'abstraction et « Tâches de raisonnement » dans l'atelier AAAI 2024 « Les grands modèles linguistiques sont-ils simplement des perroquets causaux ? » (Vancouver, BC, Canada, 2024) ; <https://openreview.net/forum?id=3rGT5OkzpC>.
- 154 L. Berglund, M. Tong, M. Kaufmann, M. Balesni, AC Stickland, T. Korbak, O. Evans, « La malédiction de l'inversion : les LLM « Formé sur « A est B », ne parvient pas à apprendre « B est A » lors de la 12e Conférence internationale sur les représentations d'apprentissage (ICLR 2024) (Vienne, Autriche, 2024) ; <https://openreview.net/forum?id=GPKTIktA0k>.
- 155 J. Geiping, A. Stein, M. Shu, K. Saifullah, Y. Wen, T. Goldstein, « Contraindre les LLM à faire et révéler (presque) n'importe quoi » dans ICLR 2024 Workshop on Secure and Trustworthy Large Language Models (SET LLM) (Vienne, Autriche, 2024) ; <https://openreview.net/forum?id=Y5inHAjMu0>.
- 156* J. Kaplan, S. McCandlish, T. Henighan, TB Brown, B. Chess, R. Child, S. Gray, A. Radford, J. Wu, D. Amodei, Lois d'échelle pour les modèles de langage neuronal, *arXiv [cs.LG]* (2020) ; <http://arxiv.org/abs/2001.08361>.
- 157* J. Hoffmann, S. Borgeaud, A. Mensch, E. Buchatskaya, T. Cai, E. Rutherford, D. de Las Casas, LA Hendricks, J. Welbl, A. Clark, T. Hennigan, E. Noland, K. Millican, G. van den Driessche, B. Damoc, A. Guy, S. Osindero, ... L. Sifre, Formation de modèles de langage de grande taille optimaux en termes de calcul, *arXiv [cs.CL]* (2022) ; <http://arxiv.org/abs/2203.15556>.
- 158* T. Henighan, J. Kaplan, M. Katz, M. Chen, C. Hesse, J. Jackson, H. Jun, TB Brown, P. Dhariwal, S. Gray, C. Hallacy, B. Mann, A. Radford, A. Ramesh, N. Ryder, DM Ziegler, J. Schulman, ... S. McCandlish, Lois d'échelle pour

- Modélisation générative autorégressive, arXiv [cs.LG] (2020) ; <http://arxiv.org/abs/2010.14701>.
- 159 X. Zhai, A. Kolesnikov, N. Houlsby, L. Beyer, « Scaling Vision Transformers » dans la conférence IEEE/CVF 2022 sur Vision par ordinateur et reconnaissance de formes (CVPR) (2022), pp. 1204–1213 ; <https://doi.org/10.1109/CVPR52688.2022.011179>.
- 160* AL Jones, Mise à l'échelle des lois d'échelle avec des jeux de société, arXiv [cs.LG] (2021) ; <http://arxiv.org/abs/2104.03113>.
- 161* Y. Bahri, E. Dyer, J. Kaplan, J. Lee, U. Sharma, Explication des lois de mise à l'échelle neuronale, arXiv [cs.LG] (2021) ; <http://arxiv.org/abs/2102.06701>.
- 162* A. Maloney, DA Roberts, J. Sully, Un modèle soluble des lois d'échelle neuronale, arXiv [cs.LG] (2022) ; <http://arxiv.org/abs/2210.16859>.
- 163 U. Sharma, J. Kaplan, Lois d'échelle de la dimension de la variété de données. *Journal of Machine Learning Research : JMLR* 23, 343–376 (2022) ; <https://dl.acm.org/doi/abs/10.5555/3586589.3586598>.
- 164 Ł. Debowski, Un modèle simpliste des lois d'échelle neuronale : processus de Santa Fe multipériodiques, arXiv [cs.IT] (2023) ; <http://arxiv.org/abs/2302.09049>.
- 165 EJ Michaud, Z. Liu, U. Girit, M. Tegmark, « Le modèle de quantification de la mise à l'échelle neuronale » dans 37e Conférence sur les systèmes de traitement de l'information neuronale (NeurIPS 2023) (La Nouvelle-Orléans, LA, États-Unis, 2023) ; <https://openreview.net/forum?id=3tbTw2ga8K>.
- 166* T. Besiroglu, E. Erdil, M. Barnett, J. You, Chinchilla Scaling : une tentative de réplication, arXiv [cs.AI] (2024) ; <http://arxiv.org/abs/2404.10102>.
- 167 T. Porian, M. Wortsman, J. Jitsev, L. Schmidt, Y. Carmon, « Résoudre les divergences dans la mise à l'échelle optimale de calcul Modèles de langage » dans le 2e atelier sur l'avancement de la formation des réseaux neuronaux : efficacité informatique, évolutivité et optimisation des ressources (WANT@ICML 2024) (2024) ; <https://openreview.net/forum?id=zhCBrgaQZO>.
- 168* T. Pearce, J. Song, Réconcilier les lois de mise à l'échelle de Kaplan et de Chinchilla, arXiv [cs.LG] (2024) ; <http://arxiv.org/abs/2406.12907>.
- 169 E. Caballero, K. Gupta, I. Rish, D. Krueger, « Lois de mise à l'échelle neuronale brisées » dans NeurIPS ML Safety Workshop (2022) ; <https://openreview.net/forum?id=BfGrFuNyhJ>.
- 170* S. Hooker, Sur les limites des seuils de calcul comme stratégie de gouvernance, arXiv [cs.AI] (2024) ; <http://arxiv.org/abs/2407.05694>.
- 171 S. Biderman, US Prashanth, L. Sutawika, H. Schoelkopf, QG Anthony, S. Purohit, E. Raff, « Mémoire émergente et prévisible dans les grands modèles linguistiques » dans 37e Conférence sur les systèmes de traitement de l'information neuronale (NeurIPS 2023) (La Nouvelle-Orléans, LA, États-Unis, 2023) ; <https://openreview.net/forum?id=lq0DvhB4Kf>.
- 172 D. Ganguli, D. Hernandez, L. Lovitt, A. Askell, Y. Bai, A. Chen, T. Conerly, N. Dassarma, D. Drain, N. Elhage, S. El Showk, S. Fort, Z. Hatfield-Dodds, T. Henighan, S. Johnston, A. Jones, N. Joseph, ... J. Clark, « Prévisibilité et surprise dans les grands modèles génératifs » dans Actes de la conférence 2022 de l'ACM sur l'équité, la responsabilité et la transparence (FAccT '22) (Association for Computing Machinery, New York, NY, États-Unis, 2022), pp. 1747–1764 ; <https://doi.org/10.1145/3531146.3533229>.
- 173* Z. Du, A. Zeng, Y. Dong, J. Tang, Comprendre les capacités émergentes des modèles linguistiques du point de vue de la perte, arXiv [cs.CL] (2024) ; <http://arxiv.org/abs/2403.15796>.
- 174 J. Wei, Y. Tay, R. Bommasani, C. Raffel, B. Zoph, S. Borgeaud, D. Yogatama, M. Bosma, D. Zhou, D. Metzler, EH Chi, T. Hashimoto, O. Vinyals, P. Liang, J. Dean, W. Fedus, Capacités émergentes des grands modèles linguistiques. *Transactions sur la recherche sur l'apprentissage automatique* (2022) ; <https://openreview.net/forum?id=yzkSU5zdwD>.
- 175 SY Gadre, G. Smyrnis, V. Shankar, S. Gururangan, M. Wortsman, R. Shao, J. Mercat, A. Fang, J. Li, S. Keh, R. Xin, M. Nezhurina, I. Vasiljevic, J. Jitsev, L. Soldaini, AG Dimakis, G. Ilharco, ... L. Schmidt, Les modèles de langage évoluent de manière fiable avec un surentraînement et des tâches en aval, arXiv [cs.CL] (2024) ; <http://arxiv.org/abs/2403.08540>.
- 176 R. Schaeffer, B. Miranda, S. Koyejo, « Les capacités émergentes des grands modèles linguistiques sont-elles un mirage ? » dans 37e Conférence sur les systèmes de traitement de l'information neuronale (NeurIPS 2023) (La Nouvelle-Orléans, LA, États-Unis, 2023) ; <https://openreview.net/forum?id=ITw9edRDID>.
- 177 Y. Ruan, CJ Maddison, T. Hashimoto, « Lois de mise à l'échelle observationnelle et prévisibilité des performances du modèle de langage » dans 38e conférence annuelle sur les systèmes de traitement de l'information neuronale (NeurIPS 2024) (2024) ; <https://openreview.net/pdf?id=On5WIN7xyD>.
- 178 TR McIntosh, T. Susnjak, T. Liu, P. Watters, MN Halgamuge, Insuffisances des repères de modèles de langage à grande échelle à l'ère de l'intelligence artificielle générative, arXiv [cs.AI] (2024) ; <http://arxiv.org/abs/2402.09880>.
- 179* V. Balachandran, J. Chen, N. Joshi, B. Nushi, H. Palangi, E. Salinas, V. Vineet, J. Woffinden-Luey, S. Yousefi, « EUREKA : évaluation et compréhension des grands modèles de fondations » (Microsoft, 2024) ;

- <https://www.microsoft.com/en-us/research/publication/eureka-evaluating-and-understanding-large-foundation-models/> .
- 180* S. Srivastava, MB Annarose, PV Anto, S. Menon, A. Sukumar, ST Adwaitha, A. Philipose, S. Prince, S. Thomas, Repères fonctionnels pour une évaluation robuste des performances de raisonnement et de l'écart de raisonnement, arXiv [cs.AI] (2024) ; <http://arxiv.org/abs/2402.19450>.
- 181 C. Deng, Y. Zhao, X. Tang, M. Gerstein, A. Cohan, Étude de la contamination des données dans les repères modernes pour les grands modèles linguistiques, arXiv [cs.CL] (2023) ; <http://arxiv.org/abs/2311.09783>.
- 182 O. Sainz, J. Campos, I. García-Ferrero, J. Etxaniz, OL de Lacalle, E. Agirre, « L'évaluation du TAL en difficulté : sur la nécessité de mesurer la contamination des données LLM pour chaque référence » dans Findings of the Association for Computational Linguistics: EMNLP 2023, H. Bouamor, J. Pino, K. Bali, éd. (Association for Computational Linguistics, Singapour, 2023), pp. 10776–10787 ; <https://doi.org/10.18653/v1/2023.findings-emnlp.722>.
- 183 Y. Cao, L. Zhou, S. Lee, L. Cabello, M. Chen, D. Hershcovich, « Évaluation de l'alignement interculturel entre ChatGPT et sociétés humaines : une étude empirique » dans les actes du 1er atelier sur les considérations interculturelles en PNL (C3NLP), S. Dev, V. Prabhakaran, D. Adelman, D. Hovy, L. Benotti, éd. (Association for Computational Linguistics, Dubrovnik, Croatie, 2023), pp. 53–67 ; <https://doi.org/10.18653/v1/2023.c3nlp-1.7>.
- 184* H. Zhou, A. Bradley, E. Littwin, N. Razin, O. Saremi, J. Susskind, S. Bengio, P. Nakkiran, What Algorithms Can Les transformateurs apprennent-ils ? Une étude sur la généralisation des longueurs, arXiv [cs.LG] (2023) ; <http://arxiv.org/abs/2310.16028>.
- 185 D. Yu, S. Kaur, A. Gupta, J. Brown-Cohen, A. Goyal, S. Arora, « SKILL-MIX : une famille flexible et extensible de Évaluations des modèles d'IA » dans la 12e Conférence internationale sur les représentations d'apprentissage (2024) ; <https://openreview.net/pdf?id=Jf5gplvqlq>.
- 186* H. Zhang, J. Da, D. Lee, V. Robinson, C. Wu, W. Song, T. Zhao, P. Raja, D. Slack, Q. Lyu, S. Hendryx, R. Kaplan, M. Lunati, S. Yue, Un examen attentif des performances du grand modèle linguistique sur l'arithmétique de l'école primaire, arXiv [cs.CL] (2024) ; <http://arxiv.org/abs/2405.00332>.
- 187* AlphaProof, équipes AlphaGeometry, l'IA atteint la norme de la médaille d'argent en résolvant les problèmes de l'Olympiade internationale de mathématiques, Google DeepMind (2024) ; <https://deepmind.google/discover/blog/ai-solves-imo-problems-at-silver-medal-level/>.
- 188 TH Trinh, Y. Wu, QV Le, H. He, T. Luong, Résolution de la géométrie olympique sans démonstration humaine. *Nature* 625, 476–482 (2024) ; <https://doi.org/10.1038/s41586-023-06747-5>.
- 189 E. Akyürek, M. Damani, L. Qiu, H. Guo, Y. Kim, J. Andreas, L'efficacité surprenante de la formation au test pour Raisonnement abstrait, arXiv [cs.AI] (2024) ; <http://arxiv.org/abs/2411.07279>.
- 190 Y. Bengio, G. Hinton, A. Yao, D. Song, P. Abbeel, T. Darrell, YN Harari, Y.-Q. Zhang, L. Xue, S. Shalev-Shwartz, G. Hadfield, J. Clune, T. Maharaj, F. Hutter, AG Baydin, S. McIlraith, Q. Gao, ... S. Mindermann, Gérer les risques extrêmes liés à l'IA dans un contexte de progrès rapide. *Science*, eadn0117 (2024) ; <https://doi.org/10.1126/science.adn0117>.
- 191 Y. LeCun, La puissance et les limites de l'apprentissage profond : dans son discours de remise de la médaille de l'IRI, Yann LeCun cartographie le développement des techniques d'apprentissage automatique et suggère ce que l'avenir pourrait réserver. *Research Technology Management* 61, 22–27 (2018) ; <https://doi.org/10.1080/08956308.2018.1516928>.
- 192 M. Mitchell, « Pourquoi l'IA est plus difficile que nous le pensons » dans Proceedings of the Genetic and Evolutionary Computation Conference (GECCO '21) (Association for Computing Machinery, New York, NY, États-Unis, 2021), p. 3 ; <https://doi.org/10.1145/3449639.3465421>.
- 193 J. Pearl, D. Mackenzie, Le livre du pourquoi : la nouvelle science de la cause et de l'effet (Penguin Books, Harlow, Angleterre, 2019) Science des pingouins ; <https://dl.acm.org/doi/10.5555/3238230>.
- 194 DC Cireşan, U. Meier, LM Gambardella, J. Schmidhuber, Réseaux neuronaux profonds, grands et simples pour les chiffres manuscrits Reconnaissance. *Neural Computation* 22, 3207–3220 (2010) ; https://doi.org/10.1162/NECO_a_00052.
- 195 T. Mikolov, M. Karafiát, L. Burget, J. Černocký, S. Khudanpur, « Modèle linguistique basé sur un réseau neuronal récurrent » dans Proc. Interspeech 2010 (ISCA, 2010), pp. <https://doi.org/10.21437/Interspeech.2010-343>.
- 196 X. Glorot, Y. Bengio, « Comprendre la difficulté de former des réseaux neuronaux à rétroaction profonde » dans Actes de la 13e Conférence internationale sur l'intelligence artificielle et les statistiques (AISTATS 2010), Yee Whye Teh, Mike Titterton, éd. (PMLR, 2010) vol. 9, pp. 249–256 ; <https://proceedings.mlr.press/v9/glorot10a.html>.
- 197 Epoch AI, Données sur les modèles d'IA notables. (2024) ; <https://epochai.org/data/notable-ai-models>.
- 198* Inflection AI, Inflection-2 (2023) ; <https://inflection.ai/inflection-2>.
- 199 C.-J. Wu, R. Raghavendra, U. Gupta, B. Acun, N. Ardalani, K. Maeng, G. Chang, F. Aga, J. Huang, C. Bai, M. Gschwind, A. Gupta, M. Ott, A. Melnikov, S. Candido, D. Brooks, G. Chauhan, ... K. Hazelwood, « IA durable : implications environnementales, défis et opportunités » dans Actes de la 5e Conférence sur l'apprentissage automatique et les systèmes (MLSys), D. Marculescu, Y. Chi, C. Wu, éd. (2022) vol. 4, pp. 795–813;

https://proceedings.mlsys.org/paper_files/paper/2022/file/462211f67c7d858f663355eff93b745e-Paper.pdf.

- 200* Y. Wu, Z. Sun, S. Li, S. Welleck, Y. Yang, Lois d'échelle d'inférence : une analyse empirique de l'optimalité de calcul
Inférence pour la résolution de problèmes avec des modèles de langage, arXiv [cs.AI] (2024) ; <http://arxiv.org/abs/2408.00724>.
- 201 S. Hao, Y. Gu, H. Ma, JJ Hong, Z. Wang, DZ Wang, Z. Hu, « Raisonner avec le modèle linguistique, c'est planifier avec le modèle mondial » dans la Conférence 2023 sur les méthodes empiriques en traitement du langage naturel (2023) ; <https://openreview.net/pdf?id=VTWWvYtF1R>.
- 202* X. Feng, Z. Wan, M. Wen, Y. Wen, W. Zhang, J. Wang, « La recherche arborescente de type Alphazero peut guider les langages de grande taille
« Décodage et formation des modèles » dans l'atelier NeurIPS 2023 Foundation Models for Decision Making (La Nouvelle-Orléans, LA, États-Unis, 2023) ; <https://openreview.net/pdf?id=PJfc4x2jXY>.
- 203* C. Li, W. Wang, J. Hu, Y. Wei, N. Zheng, H. Hu, Z. Zhang, H. Peng, Les modèles de langage courants 7B possèdent déjà de fortes capacités mathématiques, arXiv [cs.CL] (2024) ; <http://arxiv.org/abs/2403.04706>.
- 204 E. Erdil, Répartition optimale des ressources de calcul entre l'inférence et la formation. (2024) ; <https://epochai.org/blog/optimally-allocating-compute-between-inference-and-training>.
- 205 K. Chow, Y. Tang, Z. Lyu, A. Rajput, K. Ban, « Optimisation des performances dans le monde du LLM 2024 » dans Companion of the 15th ACM/ SPEC International Conference on Performance Engineering (ACM, New York, NY, États-Unis, 2024) ; <https://doi.org/10.1145/3629527.3651436>.
- 206 D. Patterson, J. Gonzalez, U. Hölzle, Q. Le, C. Liang, L.-M. Munguia, D. Rothchild, DR So, M. Texier, J. Dean, L'empreinte carbone de la formation en apprentissage automatique va plafonner, puis diminuer. *Computer* 55, 18–28 (2022) ; <https://doi.org/10.1109/MC.2022.3148714>.
- 207 D. Coyle, L. Hampton, Progrès informatiques au 21e siècle. *Telecommunications Policy* 48, 102649 (2024) ; <https://doi.org/10.1016/j.telpol.2023.102649>.
- 208 Agence internationale de l'énergie, « Électricité 2024 : analyse et prévisions jusqu'en 2026 » (AIE, 2024) ; [https://iea.blob.core.windows.net/assets/6b2fd954-2017-408e-bf08-952fdd62118a/Electricity2024-Analyse et prévisions à l'horizon 2026.pdf](https://iea.blob.core.windows.net/assets/6b2fd954-2017-408e-bf08-952fdd62118a/Electricity2024-Analyse%20et%20pr%C3%A9visions%20%C3%A0%20l'horizon%202026.pdf).
- 209 Talen Energy, Talen Energy annonce la vente d'un campus de centre de données zéro carbone (2024) ; <https://ir.talenenergy.com/news-releases/news-release-details/talen-energy-announces-sale-zero-carbon-data-center-campus>.
- 210 Pratique de l'électronique avancée, H. Bauer, O. Burkacky, P. Kenevan, S. Lingemann, K. Pototzky, B. Wiseman, « Conception et fabrication de semi-conducteurs : atteindre des capacités de pointe » (McKinsey & Company, 2020) ; https://www.mckinsey.com/industries/industrials-and-electronics/our-insights/semiconductor-design-and-manufacturing-achieving-leading-edge-capabilities#.
- 211 J. VerWey, « Pas de permis, pas de Fab : l'importance de la réforme réglementaire pour la fabrication de semi-conducteurs » (Centre pour la sécurité et les technologies émergentes, 2021) ; <https://doi.org/10.51593/20210053>.
- 212 D. Bragg, N. Caselli, JA Hochgesang, M. Huenerfauth, L. Katz-Hernandez, O. Koller, R. Kushalnagar, C. Vogler, RE Ladner, Le paysage FATE des ensembles de données d'IA en langue des signes : une perspective interdisciplinaire. *ACM Transactions on Accessible Computing* 14, 1–45 (2021) ; <https://doi.org/10.1145/3436996>.
- 213 G. Li, Z. Sun, Q. Wang, S. Wang, K. Huang, N. Zhao, Y. Di, X. Zhao, Z. Zhu, Centre de données vert de Chine Développement : Politiques et voie technologique de réduction du carbone. *Environmental Research* 231, 116248 (2023) ; <https://doi.org/10.1016/j.envres.2023.116248>.
- 214 E. Griffith, La chasse désespérée au prix le plus indispensable du boom de l'IA, *The New York Times* (2023) ; <https://www.nytimes.com/2023/08/16/technology/ai-gpu-chips-shortage.html>.
- 215 J. Sevilla, T. Besiroglu, B. Cottier, J. You, E. Roldán, P. Villalobos, E. Erdil, La mise à l'échelle de l'IA peut-elle se poursuivre jusqu'en 2030 ? (2024) ; <https://epochai.org/blog/can-ai-scaling-continue-through-2030>.
- 216 E. Erdil, « Goulots d'étranglement du mouvement des données pour la formation de modèles à grande échelle : mise à l'échelle au-delà du FLOP 1e28 » (Epoch AI, 2024) ; <https://epoch.ai/blog/data-movement-bottlenecks-scaling-past-1e28-flop>.
- 217* E. Almazrouei, H. Alobeidli, A. Alshamsi, A. Cappelli, R. Cojocar, M. Debbah, É. Goffinet, D. Hesslow, J. Launay, Q. Malartic, D. Mazzotta, B. Noune, B. Pannier, G. Penedo, La série Falcon de modèles de langage ouvert, arXiv [cs.CL] (2023) ; <http://arxiv.org/abs/2311.16867>.
- 218* T. Wei, L. Zhao, L. Zhang, B. Zhu, L. Wang, H. Yang, B. Li, C. Cheng, W. Lü, R. Hu, C. Li, L. Yang, X. Luo, X. Wu, L. Liu, W. Cheng, P. Cheng, ... Y. Zhou, Skywork : un modèle de fondation bilingue plus ouvert, arXiv [cs.CL] (2023) ; <http://arxiv.org/abs/2310.19341>.
- 219 P. Villalobos, J. Sevilla, L. Heim, T. Besiroglu, M. Hobbhahn, A. Ho, Allons-nous manquer de données ? Limites de la mise à l'échelle du LLM basée sur des données générées par l'homme, arXiv [cs.LG] (2022) ; <http://arxiv.org/abs/2211.04325>.
- 220 N. Muennighoff, A. Rush, B. Barak, T. Le Scao, N. Tazi, A. Piktus, S. Pyysalo, T. Wolf, CA Raffel, « Scaling Data-

- « Modèles de langage contraints » dans *Advances in Neural Information Processing Systems 36 (NeurIPS 2023) Conférence principale (La Nouvelle-Orléans, LA, États-Unis, 2023)* vol. 36, pp. 50358–50376 ; https://proceedings.neurips.cc/paper_files/paper/2023/hash/9d89448b63ce1e2e8dc7af72c984c196-Résumé-Conférence.html.
- 221* A. Sohn, A. Nagabandi, C. Florensa, D. Adelberg, D. Wu, H. Farooq, I. Clavera, J. Welborn, J. Chen, N. Mishra, P. Chen, P. Qian, P. Abbeel, R. Duan, V. Vijay, Y. Liu, Présentation de RFM-1 : donner aux robots des capacités de raisonnement semblables à celles des humains, *covariant (2024)* ; <https://covariant.ai/insights/introducing-rfm-1-giving-robots-human-like-reasoning-capacités/>.
- 222 H. Abdine, M. Chatzianastasis, C. Bouyioukos, M. Vazirgiannis, « Prot2Text : Fonction de la protéine multimodale Génération avec GNN et transformateurs » dans la 37e conférence sur les systèmes de traitement de l'information neuronale (NeurIPS 2023) Atelier sur les modèles génératifs profonds pour la santé (La Nouvelle-Orléans, LA, États-Unis, 2023) ; <https://openreview.net/forum?id=EJ7YNgWYFj>.
- 223 A. Radford, JW Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, G. Krueger, I. Sutskever, « Apprentissage de modèles visuels transférables à partir de la supervision du langage naturel » dans *Actes de la 38e Conférence internationale sur l'apprentissage automatique (ICML 2021) (PMLR, 2021)*, pp. 8748–8763 ; <https://proceedings.mlr.press/v139/radford21a.html>.
- 224* Communication transparente, L. Barrault, Y.-A. Chung, MC Meglioli, D. Dale, N. Dong, P.-A. Duquenne, H. Elsahar, H. Gong, K. Heffernan, J. Hoffman, C. Klaiber, P. Li, D. Licht, J. Maillard, A. Rakotoarison, KR Sadagopan, ... S. Wang, « SeamlessM4T : traduction automatique massivement multilingue et multimodale » (*Meta AI, 2023*) ; <http://arxiv.org/abs/2308.11596>.
- 225 P. Villalobos, A. Ho, J. Sevilla, T. Besiroglu, L. Heim, M. Hobbhahn, « Position : allons-nous manquer de données ? Limites de la mise à l'échelle LLM basée sur des données générées par l'homme » dans *Actes de la 41e Conférence internationale sur l'apprentissage automatique, R. Salakhutdinov, Z. Kolter, K. Heller, A. Weller, N. Oliver, J. Scarlett, F. Berkenkamp, éd. (PMLR, 2024)* vol. 235 des Actes de la recherche sur l'apprentissage automatique, pp. 49523–49544 ; <https://proceedings.mlr.press/v235/villalobos24a.html>.
- 226* L. Fan, K. Chen, D. Krishnan, D. Katabi, P. Isola, Y. Tian, Lois de mise à l'échelle des images synthétiques pour la formation de modèles... pour Maintenant, *arXiv [cs.CV] (2023)* ; <http://arxiv.org/abs/2312.04567>.
- 227 S. Fu, NY Tamir, S. Sundaram, L. Chai, R. Zhang, T. Dekel, P. Isola, « DreamSim : apprentissage de nouvelles dimensions de la similarité visuelle humaine à l'aide de données synthétiques » dans 37e Conférence sur les systèmes de traitement de l'information neuronale (NeurIPS 2023) (La Nouvelle-Orléans, LA, États-Unis, 2023) ; <https://openreview.net/forum?id=DEiNSfh1k7>.
- 228 Y. Tian, L. Fan, P. Isola, H. Chang, D. Krishnan, « StableRep : les images synthétiques issues de modèles texte-image permettent « Strong Visual Representation Learners » dans la 37e Conférence sur les systèmes de traitement de l'information neuronale (NeurIPS 2023) (La Nouvelle-Orléans, LA, États-Unis, 2023) ; <https://openreview.net/forum?id=xpjsOQtKqx>.
- 229 I. Shumailov, Z. Shumaylov, Y. Zhao, Y. Gal, N. Papernot, R. Anderson, La malédiction de la récursivité : formation sur Les données générées font oublier les modèles, *arXiv [cs.LG] (2023)* ; <http://arxiv.org/abs/2305.17493>.
- 230 G. Martínez, L. Watson, P. Reviriego, JA Hernández, M. Juarez, R. Sarkar, Combinaison de l'intelligence artificielle générative (IA) et d'Internet : vers une évolution ou une dégradation ?, *arXiv [cs.CV] (2023)* ; <http://arxiv.org/abs/2303.01255>.
- 231 R. Hataya, H. Bao, H. Arai, « Les modèles génératifs à grande échelle corrompent-ils les futurs ensembles de données ? » dans Conférence internationale IEEE/CVF 2023 sur la vision par ordinateur (ICCV) (IEEE, 2023), pp. 20498–20508 ; <https://doi.org/10.1109/iccv51070.2023.01879>.
- 232 G. Martínez, L. Watson, P. Reviriego, JA Hernández, M. Juarez, R. Sarkar, « Vers une compréhension de l'interaction entre l'intelligence artificielle générative et Internet » dans *Lecture Notes in Computer Science (Springer Nature Suisse, Cham, 2024)* vol. 14523 de Lecture notes in computer science, pp. 59–73 ; https://doi.org/10.1007/978-3-031-57963-9_5.
- 233 Y. Guo, G. Shang, M. Vazirgiannis, C. Clavel, Le curieux déclin de la diversité linguistique : formation de modèles linguistiques sur du texte synthétique, *arXiv [cs.CL] (2023)* ; <http://arxiv.org/abs/2311.09807>.
- 234* M. Bohacek, H. Farid, Effondrement des modèles d'IA générative formés de manière népotique, *arXiv [cs.AI] (2023)* ; <http://arxiv.org/abs/2311.12202>.
- 235 S. Alemohammad, J. Casco-Rodriguez, L. Luzi, Al Humayun, H. Babaei, D. LeJeune, A. Siahkoohi, R. Baraniuk, « Les modèles génératifs autoconsomants deviennent fous » dans la 12e Conférence internationale sur les représentations d'apprentissage (ICLR 2024) (Vienne, Autriche, 2023) ; <https://openreview.net/forum?id=ShjMHfmPs0>.
- 236 Q. Bertrand, J. Bose, A. Duplessis, M. Jiralerspong, G. Gidel, « Sur la stabilité du recyclage itératif des modèles génératifs sur leurs propres données » dans 12e Conférence internationale sur les représentations d'apprentissage (2024) ; <https://openreview.net/forum?id=JORAfH2xFd>.

- 237* E. Dohmatob, Y. Feng, P. Yang, F. Charton, J. Kempe, A Tale of Tails : l'effondrement du modèle comme changement des lois d'échelle, arXiv [cs.LG] (2024) ; <http://arxiv.org/abs/2402.07043>.
- 238 R. He, S. Sun, X. Yu, C. Xue, W. Zhang, P. Torr, S. Bai, X. Qi, « Les données synthétiques issues de modèles génératifs sont-elles prêtes pour Reconnaissance d'images ? » dans la 11e Conférence internationale sur les représentations d'apprentissage (ICLR 2023) (Kigali, Rwanda, 2022) ; <https://openreview.net/pdf?id=nUmCcZ5RKF>.
- 239* V. Boutin, L. Singhal, X. Thomas, T. Serre, « Diversité contre reconnaissabilité : généralisation à la manière humaine en one-shot « Modèles génératifs » dans Advances in Neural Information Processing Systems (NeurIPS 2022) (La Nouvelle-Orléans, LA, États-Unis, 2022) ; <https://openreview.net/pdf?id=DVfZKXSFw5m>.
- 240 V. Boutin, T. Fel, L. Singhal, R. Mukherji, A. Nagaraj, J. Colin, T. Serre, « Les modèles de diffusion en tant qu'artistes : comblons-nous l'écart entre les humains et les machines ? » dans Actes de la 40e Conférence internationale sur l'apprentissage automatique (PMLR, 2023), pp. 2953–3002 ; <https://proceedings.mlr.press/v202/boutin23a.html>.
- 241 J. Shipard, A. Wiliem, KN Thanh, W. Xiang, C. Fookes, « La diversité est absolument nécessaire : améliorer la classification Zero-Shot indépendante du modèle via la diffusion stable » dans Ateliers de la conférence IEEE/CVF 2023 sur la vision par ordinateur et la reconnaissance de formes (CVPRW) (IEEE, 2023), pp. 769–778 ; <https://doi.org/10.1109/cvprw59228.2023.00084>.
- 242* A. Setlur, S. Garg, X. Geng, N. Garg, V. Smith, A. Kumar, RL sur Les données synthétiques incorrectes multiplient par huit l'efficacité du raisonnement mathématique LLM, arXiv [cs.LG] (2024) ; <http://arxiv.org/abs/2406.14532>.
- 243 P. Haluptzok, M. Bowers, AT Kalai, « Les modèles de langage peuvent apprendre à mieux programmer » dans Deep Atelier d'apprentissage par renforcement NeurIPS 2022 (2022) ; https://openreview.net/forum?id=_5BZwkZRFc9.
- 244* B. Liu, S. Bubeck, R. Eldan, J. Kulkarni, Y. Li, A. Nguyen, R. Ward, Y. Zhang, TinyGSM : atteindre > 80 % sur GSM8k avec Petits modèles de langage, arXiv [cs.LG] (2023) ; <http://arxiv.org/abs/2312.09241>.
- 245* D. Hernandez, TB Brown, Mesure de l'efficacité algorithmique des réseaux neuronaux, arXiv [cs.LG] (2020) ; <http://arxiv.org/abs/2005.04305>.
- 246 A. Ho, T. Besiroglu, E. Erdil, D. Owen, R. Rahman, ZC Guo, D. Atkinson, N. Thompson, J. Sevilla, « Progrès algorithmiques dans les modèles linguistiques » (Epoch AI, 2024) ; <http://arxiv.org/abs/2403.05812>.
- 247 FE Dorner, Mesure des progrès dans l'efficacité des échantillons d'apprentissage par renforcement profond, arXiv [cs.LG] (2021) ; <http://arxiv.org/abs/2102.04881>.
- 248* Y. Ding, LL Zhang, C. Zhang, Y. Xu, N. Shang, J. Xu, F. Yang, M. Yang, LongRoPE : Extension de la fenêtre contextuelle LLM au-delà de 2 millions de jetons, arXiv [cs.CL] (2024) ; <http://arxiv.org/abs/2402.13753>.
- 249 A. Fawzi, M. Balog, A. Huang, T. Hubert, B. Romera-Paredes, M. Barekatin, A. Novikov, FJ R Ruiz, J. Schrittwieser, G. Swirszcz, D. Silver, D. Hassabis, P. Kohli, Découverte d'algorithmes de multiplication de matrices plus rapides grâce à l'apprentissage par renforcement. Nature 610, 47–53 (2022) ; <https://doi.org/10.1038/s41586-022-05172-4>.
- 250 A. Haj-Ali, NK Ahmed, T. Willke, YS Shao, K. Asanovic, I. Stoica, « NeuroVectorizer : vectorisation de bout en bout avec apprentissage par renforcement profond » dans Actes du 18e Symposium international ACM/IEEE sur la génération et l'optimisation de code (CGO 2020) (Association for Computing Machinery, New York, NY, États-Unis, 2020), pp. 242–255 ; <https://doi.org/10.1145/3368826.3377928>.
- 251 A. Goldie, A. Mirhoseini, M. Yazgan, JW Jiang, E. Songhori, S. Wang, Y.-J. Lee, E. Johnson, O. Pathak, A. Nova, J. Pak, A. Tong, K. Srinivasa, W. Hang, E. Tuncer, QV Le, J. Laudon, ... J. Dean, Addendum : Une méthodologie de placement de graphes pour une conception rapide de puces. Nature 634, E10–E11 (2024) ; <https://doi.org/10.1038/s41586-024-08032-5>.
- 252 X. Li, P. Yu, C. Zhou, T. Schick, O. Levy, L. Zettlemoyer, JE Weston, M. Lewis, « Auto-alignement avec l'instruction « Rétrotraduction » dans la 12e Conférence internationale sur les représentations de l'apprentissage (ICLR 2024) (Vienne, Autriche, 2023) ; <https://openreview.net/forum?id=1oijHJBRsT>.
- 253 S. Liu, Z. Lin, S. Yu, R. Lee, T. Ling, D. Pathak, D. Ramanan, Modèles de langage comme optimiseurs de boîte noire pour la vision Modèles de langage, arXiv [cs.CL] (2023) ; <http://arxiv.org/abs/2309.05950>.
- 254 R. Pryzant, D. Iyer, J. Li, Y. Lee, C. Zhu, M. Zeng, « Optimisation automatique des invites avec « descente de gradient » et « Beam Search » dans les actes de la conférence 2023 sur les méthodes empiriques en traitement du langage naturel (EMNLP 2023), H. Bouamor, J. Pino, K. Bali, éd. (Association for Computational Linguistics, Singapour, 2023), pp. 7957–7968 ; <https://doi.org/10.18653/v1/2023.emnlp-main.494>.
- 255 S. Zhang, C. Gong, L. Wu, X. Liu, M. Zhou, AutoML-GPT : apprentissage automatique avec GPT, arXiv [cs.CL] (2023) ; <http://arxiv.org/abs/2305.02499>.
- 256* Y. Bai, S. Kadavath, S. Kundu, A. Askell, J. Kernion, A. Jones, A. Chen, A. Goldie, A. Mirhoseini, C. McKinnon, C. Chen, C. Olsson, C. Olah, D. Hernandez, D. Drain, D. Ganguli, D. Li, ... J. Kaplan, IA constitutionnelle : innocuité des commentaires sur l'IA, arXiv [cs.CL] (2022) ; <http://arxiv.org/abs/2212.08073>.
- 257* N. Sachdeva, B. Coleman, W.-C. Kang, J. Ni, L. Hong, EH Chi, J. Caverlee, J. McAuley, DZ Cheng, Comment s'entraîner LLMs efficaces en termes de données, arXiv [cs.LG] (2024) ; <http://arxiv.org/abs/2402.09668>.

- 258* S. Kumar, T. Ghosal, V. Goyal, A. Ekbal, Les grands modèles linguistiques peuvent-ils ouvrir la voie à de nouvelles idées de recherche scientifique ?, arXiv [cs.CL] (2024) ; <http://arxiv.org/abs/2409.06185>.
- 259 H. Wijk, T. Lin, J. Becker, S. Jawhar, N. Parikh, T. Broadley, L. Chan, M. Chen, J. Clymer, J. Dhyani, E. Elicheva, K. Garcia, B. Goodrich, N. Jurkovic, M. Kinniment, A. Lajko, S. Nix, ... E. Barnes, RE-Bench : Évaluation des capacités de R&D en IA de pointe des agents de modèles linguistiques par rapport aux experts humains, arXiv [cs.LG] (2024) ; <http://arxiv.org/abs/2411.15114>.
- 260 D. Owen, « Entretiens avec des chercheurs en IA sur l'automatisation de la R&D en IA » (Epoch AI, 2024) ; <https://epoch.ai/blog/interviewing-ai-researchers-on-automation-of-ai-rnd>.
- 261* E. Erdil, J. Sevilla, Tendances des lois de puissance dans le speedrunning et l'apprentissage automatique, arXiv [cs.LG] (2023) ; <http://arxiv.org/abs/2304.10004>.
- 262* J. Droppo, O. Elibol, Lois d'échelle pour les modèles acoustiques, arXiv [eess.AS] (2021) ; <http://arxiv.org/abs/2106.09488>.
- 263 S. Hooker, La loterie du matériel. Communications de l'ACM 64, 58–65 (2021) ; <https://doi.org/10.1145/3467017>.
- 264* Q. Anthony, J. Hatef, D. Narayanan, S. Biderman, S. Bekman, J. Yin, A. Shafi, H. Subramoni, D. Panda, Plaidoyer pour la co-conception d'architectures de modèles avec du matériel, arXiv [cs.DC] (2024) ; <http://arxiv.org/abs/2401.14489>.
- 265* F. Mince, D. Dinh, J. Kgomo, N. Thompson, S. Hooker, La grande illusion : le mythe de la portabilité des logiciels et Implications pour les progrès du ML, arXiv [cs.SE] (2023) ; <http://arxiv.org/abs/2309.07181>.
- 266* L'équipe Scale, soumettez vos questions les plus difficiles pour le dernier examen de l'humanité, échelle (2024) ; <https://scale.com/blog/humanitys-last-exam>.
- 267 Prix ARC, Prix ARC, Prix ARC (2024) ; <https://arcprize.org/>.
- 268 Département des sciences, de l'innovation et de la technologie, « Approche de l'AI Safety Institute en matière d'évaluations » (GOV.UK, 2024) ; <https://www.gov.uk/government/publications/ai-safety-institute-approach-to-evaluations/ai-safety-institute-approach-to-evaluations> .
- 269 Metr, Une mise à jour sur nos évaluations de capacité générale, METR (2024) ; <https://metr.org/blog/2024-08-06-mise-à-jour-des-évaluations/>.
- 270 G. Sastry, L. Heim, H. Belfield, M. Anderljung, M. Brundage, J. Hazell, C. O'Keefe, GK Hadfield, R. Ngo, K. Pilz, G. Gor, E. Bluemke, S. Shoker, J. Egan, RF Trager, S. Avin, A. Weller, ... D. Coyle, La puissance de calcul et la gouvernance de l'intelligence artificielle, arXiv [cs.CY] (2024) ; <http://arxiv.org/abs/2402.08797>.
- 271 D. Citron, R. Chesney, Deep Fakes : un défi imminent pour la vie privée, la démocratie et la sécurité nationale. California Law Review 107, 1753 (2019) ; https://scholarship.law.bu.edu/faculty_scholarship/640.
- 272 Nations Unies, Déclaration universelle des droits de l'homme (1948) ; <https://www.un.org/fr/a-propos-de-nous/declaration-universelle-des-droits-de-l-homme> .
- 273 V. Ciancaglini, C. Gibson, D. Sancho, O. McCarthy, M. Eira, P. Amann, A. Klayn, « Utilisations malveillantes et abus de l'intelligence artificielle » (Agence de l'Union européenne pour la coopération des services répressifs, 2020) ; https://documents.trendmicro.com/assets/white_papers/wp-malicious-uses-and-abuses-of-artificial-intelligence.pdf .
- 274 PV Falade, Décoder le paysage des menaces : ChatGPT, FraudGPT et WormGPT dans les attaques d'ingénierie sociale. Journal international de recherche scientifique en informatique, ingénierie et technologies de l'information 9, 185–198 (2023) ; <https://doi.org/10.32628/CSEIT2390533>.
- 275 J. Bateman, « Deepfakes et médias synthétiques dans le système financier : évaluation des scénarios de menace » (Carnegie Fonds pour la paix internationale, 2020) ; <https://carnegieendowment.org/research/2020/07/deepfakes-and-synthetic-media-in-the-financial-system-assessing-threat-scenarios?lang=fr> .
- 276 Federal Bureau of Investigation des États-Unis, numéro d'alerte I-060523-PSA : acteurs malveillants manipulant des photos et Vidéos pour créer du contenu explicite et des stratagèmes de sextorsion (2023) ; <https://www.ic3.gov/PSA/2023/psa230605>.
- 277 A. Kaur, A. Noori Hoshiyar, V. Saikrishna, S. Firmin, F. Xia, Détection vidéo Deepfake : défis et opportunités. Revue d'intelligence artificielle 57, 1–47 (2024) ; <https://doi.org/10.1007/s10462-024-10810-6>.
- 278 R. Umbach, N. Henry, G. Beard, C. Berryessa, Imagerie intime synthétique non consensuelle : prévalence, attitudes, et connaissances dans 10 pays, arXiv [cs.CY] (2024) ; <http://arxiv.org/abs/2402.01721>.
- 279 Mo Kugler, C. Pace, Deepfake Privacy: Attitudes and Regulation. Revue de droit de l'université Northwestern 116, 611–680 (2021) ; <https://scholarlycommons.law.northwestern.edu/nulr/vol116/iss3/1>.
- 280 M. Viola, C. Voto, Conçus pour abuser ? Les deepfakes et la diffusion non consensuelle d'images intimes. Synthèse 201, 30 (2023) ; <https://doi.org/10.1007/s11229-022-04012-2>.
- 281 S. Maddocks, « Un complot pornographique deepfake destiné à me faire taire » : exploration des continuités entre la pornographie et

- Faux profonds « politiques ». *Porn Studies* 7, 415–423 (2020) ; <https://doi.org/10.1080/23268743.2020.1757499>.
- 282 H. Ajder, G. Patrini, F. Cavalli, L. Cullen, « L'état des Deepfakes : paysage, menaces et impact » (Deeprace, 2019); https://regmedia.co.uk/2019/10/08/deepfake_report.pdf.
- 283 J. Laffier, A. Rehman, Deepfakes et préjudices causés aux femmes. *Journal of Digital Life and Learning* 3, 1–21 (2023) ; <https://doi.org/10.51357/jdll.v3i1.218>.
- 284* T. Sippy, F. Enock, J. Bright, HZ Margetts, Derrière le Deepfake : 8 % créent ; 90 % s'inquiètent. Enquête sur l'exposition du public aux Deepfakes et les perceptions de ces derniers au Royaume-Uni, arXiv [cs.CY] (2024) ; <http://arxiv.org/abs/2407.05529>.
- 285 D. Thiel, « Identifier et éliminer le CSAM dans les données et modèles de formation générative ML » (Stanford Digital Dépôt, 2023) ; <https://purl.stanford.edu/kh752sm9123>.
- 286 Ofcom, Une plongée en profondeur dans les deepfakes qui dégradent, fraudent et désinforment (2024) ; <https://www.ofcom.org.uk/online-safety/illegal-and-harmful-content/deepfakes-demean-defraud-disinform/>.
- 287 S. Dunn, Définitions juridiques des images intimes à l'ère des deepfakes sexuels et de l'IA générative, Social Science Research Network (2024) ; <https://papers.ssrn.com/abstract=4813941>.
- 288 Y. Mirsky, W. Lee, La création et la détection de deepfakes : une enquête, arXiv [cs.CV] (2020) ; <http://arxiv.org/abs/2004.11138>.
- 289 A. Lewis, P. Vu, R. Duch, A. Chowdhury, Les avertissements de contenu aident-ils les gens à repérer un deepfake ? Preuves issues de deux expériences (2022) ; <https://royalsocietypublishing.org/doi/10.1098/rsos.220301>.
- 290 A. Qureshi, D. Megias, M. Kuribayashi, « Détection de vidéos deepfake à l'aide du filigrane numérique » en 2021 en Asie-Sommet annuel et conférence de l'Association du traitement des signaux et de l'information du Pacifique (APSIPA ASC) (2021), pp. 1786–1793 ; <http://www.apsipa.org/proceedings/2021/pdfs/0001786.pdf>.
- 291 L. Tang, Q. Ye, H. Hu, Q. Xue, Y. Xiao, J. Li, DeepMark : un cadre évolutif et robuste pour la détection de vidéos DeepFake. *ACM Transactions on Privacy and Security* 27, 1–26 (2024) ; <https://doi.org/10.1145/3629976>.
- 292 L.-Y. Hsu, Détection de deepfakes assistée par IA à l'aide du filigranage d'image adaptatif à l'aveugle. *Journal of Visual Communication et représentation d'image* 100, 104094 (2024) ; <https://doi.org/10.1016/j.jvcir.2024.104094>.
- 293 Y. Zhao, B. Liu, M. Ding, B. Liu, T. Zhu, X. Yu, « Défense proactive contre les deepfakes via le filigranage d'identité » dans Conférence d'hiver IEEE/CVF 2023 sur les applications de la vision par ordinateur (WACV) (2023), pp. 4591–4600 ; <https://doi.org/10.1109/WACV56688.2023.00458>.
- 294* S. Goyal, P. Kohli, Identification des images générées par l'IA avec SynthID, Google DeepMind (2023) ; <https://deepmind.google/discover/blog/identifying-ai-generated-images-with-synthid/>.
- 295 AJ Patil, R. Shelke, Un filigranage audio numérique efficace utilisant un réseau neuronal convolutionnel profond avec un algorithme d'optimisation de l'emplacement de recherche pour améliorer la robustesse et l'imperceptibilité. *High-Confidence Computing* 3, 100153 (2023) ; <https://doi.org/10.1016/j.hcc.2023.100153>.
- 296 MS Uddin, Ohidujjaman, M. Hasan, T. Shimamura, Filigranage audio : une revue complète. *Revue internationale d'informatique avancée et d'applications* 15 (2024) ; <https://doi.org/10.14569/IJACSA.2024.01505141>.
- 297 S. Abdelnabi, M. Fritz, « Transformateur de filigrane contradictoire : vers le traçage de la provenance du texte avec les données « Cacher » dans le Symposium IEEE sur la sécurité et la confidentialité (2021), pp. 121–140 ; <https://doi.org/10.1109/SP40001.2021.00083>.
- 298* X. Zhao, K. Zhang, Z. Su, S. Vasani, I. Grishchenko, C. Kruegel, G. Vigna, Y.-X. Wang, L. Li, Filigranes d'images invisibles Sont prouvablement amovibles à l'aide de l'IA générative, arXiv [cs.CR] (2023) ; <http://arxiv.org/abs/2306.01953>.
- 299 M. Saber, VS Sadasivan, K. Rezaei, A. Kumar, A. Chegini, W. Wang, S. Feizi, « Robustesse des détecteurs d'images IA : Limites fondamentales et attaques pratiques » dans la 12e Conférence internationale sur les représentations d'apprentissage (2023) ; <https://openreview.net/pdf?id=dLoAdIKENc>.
- 300 G. Björkstén, « Identifier le contenu de l'IA générative : quand et comment le filigrane peut contribuer à faire respecter les droits de l'homme » (accessnow, 2023) ; <https://www.accessnow.org/wp-content/uploads/2023/09/Identifying-generative-AI-content-when-and-how-watermarking-can-help-uphold-human-rights.pdf>.
- 301 D. Cooke, A. Edwards, S. Barkoff, K. Kelly, Aussi bon qu'un tirage au sort : détection humaine d'images, de vidéos, d'audio et de stimuli audiovisuels générés par l'IA, arXiv [cs.HC] (2024) ; <http://arxiv.org/abs/2403.16760>.
- 302 M. Jakesch, JT Hancock, M. Naaman, Les heuristiques humaines pour le langage généré par l'IA sont défectueuses. *Actes de l'Académie nationale des sciences des États-Unis d'Amérique* 120, e2208839120 (2023) ; <https://doi.org/10.1073/pnas.2208839120>.
- 303 G. Spitale, N. Biller-Andorno, F. Germani, Le modèle d'IA GPT-3 (dés)informe mieux que les humains. *Progrès scientifiques*

- 9, eadh1850 (2023); <https://doi.org/10.1126/sciadv.adh1850>.
- 304 S. Kreps, RM McCain, M. Brundage, Toutes les nouvelles qui méritent d'être inventées : le texte généré par l'IA comme outil de désinformation médiatique. *Journal of Experimental Political Science* 9, 104–117 (2022) ; <https://doi.org/10.1017/xps.2020.37>.
- 305 NC Köbis, B. Doležalová, I. Soraperra, Fooled Twice : les gens ne peuvent pas détecter les deepfakes mais pensent qu'ils le peuvent. *iScience* 24 (2021); <https://doi.org/10.1016/j.isci.2021.103364>.
- 306 K.-C. Yang, F. Menczer, Anatomie d'un botnet social malveillant alimenté par l'IA, arXiv [cs.CY] (2023) ; <http://arxiv.org/abs/2307.16336>.
- 307 R. Raman, V. Kumar Nair, P. Nedungadi, A. Kumar Sahu, R. Kowalski, S. Ramanathan, K. Achuthan, Fake News Tendances de la recherche, liens avec l'intelligence artificielle générative et les objectifs de développement durable. *Heliyon* 10, e24727 (2024) ; <https://doi.org/10.1016/j.heliyon.2024.e24727>.
- 308* M. Musser, Une analyse des coûts des modèles de langage génératifs et des opérations d'influence, arXiv [cs.CY] (2023) ; <http://arxiv.org/abs/2308.03740>.
- 309 H. Bai, JG Voelkel, JC Eichstaedt, R. Willer, L'intelligence artificielle peut persuader les humains sur les questions politiques (2023) ; <https://doi.org/10.31219/osf.io/stakv>.
- 310 K. Hackenburg, L. Ibrahim, BM Tappin, M. Tsakiris, Comparaison de la force de persuasion des grands modèles de langage de jeu de rôle et des experts humains sur les questions politiques polarisées des États-Unis (2023) ; <https://doi.org/10.31219/osf.io/ey8db>.
- 311 JA Goldstein, J. Chao, S. Grossman, A. Stamos, M. Tomz, Dans quelle mesure la propagande générée par l'IA est-elle persuasive ? *PNAS Nexus* 3, gae034 (2024) ; <https://doi.org/10.1093/pnasnexus/pgae034>.
- 312 SC Matz, JD Teeny, SS Vaid, H. Peters, GM Harari, M. Cerf, Le potentiel de l'IA générative pour la persuasion personnalisée à grande échelle. *Scientific Reports* 14, 4692 (2024) ; <https://doi.org/10.1038/s41598-024-53755-0>.
- 313* AR Williams, L. Burke-Moore, RS-Y. Chan, FE Enock, F. Nanni, T. Sippy, Y.-L. Chung, E. Gabasova, K. Hackenburg, J. Bright, Les grands modèles linguistiques peuvent générer systématiquement du contenu de haute qualité pour les opérations de désinformation électorale, arXiv [cs.CY] (2024) ; <http://arxiv.org/abs/2408.06731>.
- 314 TH Costello, G. Pennycook, DG Rand, Réduire durablement les croyances conspirationnistes grâce au dialogue avec l'IA. *Science* (New York, NY) 385, eadq1814 (2024) ; <https://doi.org/10.1126/science.adq1814>.
- 315 F. Salvi, MH Ribeiro, R. Gallotti, R. West, Sur la persuasion conversationnelle des grands modèles linguistiques : un essai contrôlé randomisé, arXiv [cs.CY] (2024) ; <http://arxiv.org/abs/2403.14380>.
- 316* I. Gabriel, A. Manzini, G. Keeling, LA Hendricks, V. Rieser, H. Iqbal, N. Tomašev, I. Ktena, Z. Kenton, M. Rodriguez, S. El-Sayed, S. Brown, C. Akbulut, A. Trask, E. Hughes, A. Stevie Bergman, R. Shelby, ... J. Manyika, « L'éthique des assistants IA avancés » (Google DeepMind, 2024) ; <http://arxiv.org/abs/2404.16244>.
- 317 PS Park, S. Goldstein, A. O'Gara, M. Chen, D. Hendrycks, La tromperie de l'IA : une étude des exemples, des risques et du potentiel Solutions. *Modèles* 5 (2024) ; <https://doi.org/10.1016/j.patter.2024.100988>.
- 318* M. Phuong, M. Aitchison, E. Catt, S. Cogan, A. Kaskasoli, V. Krakovna, D. Lindner, M. Rahtz, Y. Assael, S. Hodkinson, H. Howard, T. Lieberum, R. Kumar, MA Raad, A. Webson, L. Ho, S. Lin, ... T. Shevlane, « Évaluation des modèles de frontière pour les capacités dangereuses » (Google Deepmind, 2024) ; <https://doi.org/10.48550/arXiv.2403.13793>.
- 319 M. Burtell, T. Woodside, Influence artificielle : une analyse de la persuasion basée sur l'IA, arXiv [cs.CY] (2023) ; <http://arxiv.org/abs/2303.08721>.
- 320 F. Miró-Llinares, JC Aguerri, Désinformation sur les fausses nouvelles : une revue critique systématique des études empiriques sur le phénomène et son statut de « menace ». *Revue européenne de criminologie* 20, 356–374 (2023) ; <https://doi.org/10.1177/1477370821994059>.
- 321 G. Pennycook, DG Rand, Combattre la désinformation sur les médias sociaux en utilisant des jugements participatifs sur la qualité des sources d'information. *Actes de l'Académie nationale des sciences des États-Unis d'Amérique* 116, 2521–2526 (2019) ; <https://doi.org/10.1073/pnas.1806781116>.
- 322 Z. Epstein, N. Sirlin, A. Arechar, G. Pennycook, D. Rand, Le contexte des médias sociaux interfère avec le discernement de la vérité. *Science Advances* 9, eabo6169 (2023) ; <https://doi.org/10.1126/sciadv.abo6169>.
- 323 G. Pennycook, Z. Epstein, M. Mosleh, AA Arechar, D. Eckles, DG Rand, Porter une attention accrue à l'exacitude peut réduire la désinformation en ligne. *Nature* 592, 590–595 (2021) ; <https://doi.org/10.1038/s41586-021-03344-2>.
- 324 Pew Research Center, Une majorité d'Américains sont très préoccupés par le fait que l'IA sera utilisée pour créer de fausses informations à propos des candidats de 2024 (2024) ; https://www.pewresearch.org/short-reads/2024/09/19/concern-over-the-impact-of-ai-on-2024-presidential-campaign/sr_24-09-10_electionandai_01/.
- 325 S. Kapoor, A. Narayanan, « Comment se préparer au déluge d'IA générative sur les médias sociaux : une analyse fondée des défis et des opportunités » (Knight First Amendment Institute de l'Université Columbia, 2023) ;

- <https://s3.amazonaws.com/kfai-documents/documents/a566f4ded5/Comment-se-préparer-au-déluge-d'IA-généralive-sur-les-médias-sociaux.pdf> .
- 326 M. Hameleers, Manipulation bon marché contre manipulation profonde : les effets des faux bon marché contre les faux profonds dans un contexte politique. *Revue internationale de recherche sur l'opinion publique* 36 (2024) ; <https://doi.org/10.1093/ijpor/edae004>.
- 327 S. Vosoughi, D. Roy, S. Aral, La propagation des nouvelles vraies et fausses en ligne. *Science* 359, 1146–1151 (2018) ; <https://doi.org/10.1126/science.aap9559>.
- 328 K. Clayton, S. Blair, J. A. Busam, S. Forstner, J. Glance, G. Green, A. Kawata, A. Kovvuri, J. Martin, E. Morgan, M. Sandhu, R. Sang, R. Scholz-Bright, A. T. Welch, A. G. Wolff, A. Zhou, B. Nyhan, De vraies solutions pour les fausses nouvelles ? Mesure de l'efficacité des avertissements généraux et des balises de vérification des faits pour réduire la croyance aux fausses histoires sur les réseaux sociaux. *Political Behavior* 42, 1073–1095 (2020) ; <https://doi.org/10.1007/s11109-019-09533-0>.
- 329 E. Hoes, B. Aitken, J. Zhang, T. Gackowski, M. Wojcieszak, Les interventions de désinformation importantes réduisent les perceptions erronées mais augmentent le scepticisme, *PsyArXiv* (2023) ; <https://doi.org/10.31234/osf.io/zmpdu>.
- 330 A. Bashardoust, S. Feuerriegel, YR Shrestha, Comparaison de la volonté de partager entre les créations humaines et celles issues de l'IA Fausses nouvelles générées. *Actes de l'ACM sur l'interaction homme-machine* 8, 1–21 (2024) ; <https://doi.org/10.1145/3687028>.
- 331 A. Kumar, JW Taylor, Importance des fonctionnalités à l'ère de l'IA explicable : étude de cas sur la détection de fausses nouvelles et de désinformation via un cadre multimodal. *Revue européenne de recherche opérationnelle* 317, 401–413 (2024) ; <https://doi.org/10.1016/j.ejor.2023.10.003>.
- 332 SS Ghosal, S. Chakraborty, J. Geiping, F. Huang, D. Manocha, A. Bedi, Une enquête sur les possibilités et les impossibilités de la détection de texte générée par l'IA. *Transactions on Machine Learning Research* (2023) ; <https://openreview.net/pdf?id=AXtFeYjboj>.
- 333 VS Sadasivan, A. Kumar, S. Balasubramanian, W. Wang, S. Feizi, Le texte généré par l'IA peut-il être détecté de manière fiable ?, *arXiv [cs.CL]* (2023) ; <http://arxiv.org/abs/2303.11156>.
- 334 S. Gehrmann, H. Strobelt, A. Rush, « GLTR : détection statistique et visualisation du texte généré » dans *Actes de la 57e réunion annuelle de l'Association for Computational Linguistics : System Demonstrations*, MR Costa-jussà, E. Alfonseca, éd. (Association for Computational Linguistics, Florence, Italie, 2019), pp. 111–116 ; <https://doi.org/10.18653/v1/P19-3019>.
- 335 L. Fröhling, A. Zubiaga, Détection basée sur les caractéristiques des modèles de langage automatisés : s'attaquer à GPT-2, GPT-3 et Grover. *PeerJ. Informatique* 7, e443 (2021) ; <https://doi.org/10.7717/peerj-cs.443>.
- 336 J. Luo, G. Nan, D. Li, Y. Tan, Détection des faux avis générés par l'IA. (2023) ; <https://doi.org/10.2139/ssrn.4610727>.
- 337 T. Berber Sardinha, Textes générés par l'IA et textes rédigés par des humains : une comparaison multidimensionnelle. *Corpus appliqué Linguistique* 4, 100083 (2024) ; <https://doi.org/10.1016/j.acorp.2023.100083>.
- 338 DM Markowitz, JT Hancock, JN Bailenson, Marqueurs linguistiques de la communication d'IA intrinsèquement fausse et de la communication humaine intentionnellement fausse : preuves tirées des critiques d'hôtels. *Journal of Language and Social Psychology* 43, 63–82 (2024) ; <https://doi.org/10.1177/0261927X231200201>.
- 339 Y. Xie, A. Rawal, Y. Cen, D. Zhao, SK Narang, S. Sushmita, MUGC : Généré par la machine versus généré par l'utilisateur Détection de contenu, *arXiv [cs.CL]* (2024) ; <http://arxiv.org/abs/2403.19725>.
- 340 J. Su, TY Zhuo, J. Mansurov, D. Wang, P. Nakov, Les détecteurs de fausses nouvelles sont biaisés par rapport aux textes générés par de grands modèles linguistiques, *arXiv [cs.CL]* (2023) ; <http://arxiv.org/abs/2309.08674>.
- 341 W. Liang, M. Yuksekgonul, Y. Mao, E. Wu, J. Zou, « Les détecteurs GPT sont biaisés contre les auteurs non anglophones » dans *Atelier ICLR 2023 sur les modèles d'apprentissage automatique à grande échelle fiables et dignes de confiance* (2023) ; <https://openreview.net/pdf?id=SPuX8tKKIQ>.
- 342 A. Uchendu, J. Lee, H. Shen, T. Le, T.-H. 'kenneth Huang, D. Lee, La collaboration humaine améliore-t-elle la précision Comment identifier les textes Deepfake générés par LLM ?, *arXiv [cs.CL]* (2023) ; <http://arxiv.org/abs/2304.01002>.
- 343 MK Land, Contre la censure privatisée : propositions pour une délégation responsable. *Virginia Journal of International Law* 60, 363 (2019) ; https://papers.ssrn.com/sol3/papers.cfm?abstract_id=3442184.
- 344 R. Gorwa, R. Binns, C. Katzenbach, Modération algorithmique du contenu : défis techniques et politiques dans l'automatisation de la gouvernance des plateformes. *Big Data & Society* 7, 205395171989794 (2020) ; <https://doi.org/10.1177/2053951719897945>.
- 345 J. Turner, *Robot Rules* (Springer International Publishing, Cham, Suisse, éd. 1, 2018) ; <https://doi.org/10.1007/978-3-319-96235-1>.
- 346 N. Bontridder, Y. Poulet, Le rôle de l'intelligence artificielle dans la désinformation. *Data & Policy* 3, e32 (2021) ; <https://doi.org/10.1017/dap.2021.20>.

- 347 TC Helmus, Intelligence artificielle, Deepfakes et désinformation : une introduction (RAND Corporation, Santa Monica, CA, 2022) ; <https://doi.org/10.7249/PEA1043-1>.
- 348 S. Metta, I. Chang, J. Parker, MP Roman, AF Ehuan, Generative AI in Cybersecurity, arXiv [cs.CR] (2024) ; <http://arxiv.org/abs/2405.01674>.
- 349 Centre national de cybersécurité (NCSC), « L'impact à court terme de l'IA sur la cybermenace » (GOV.UK, 2024) ; <https://www.ncsc.gov.uk/report/impact-of-ai-on-cyber-threat>.
- 350 British Library, « Tirer les leçons de la cyberattaque : revue des incidents cybernétiques de la British Library » (British Library, 2024) ; <https://www.bl.uk/home/british-library-cyber-incident-review-8-march-2024.pdf>.
- 351* Microsoft Threat Intelligence, Garder une longueur d'avance sur les acteurs de la menace à l'ère de l'IA, Microsoft Security Blog (2024) ; <https://www.microsoft.com/en-us/security/blog/2024/02/14/staying-ahead-of-threat-actors-in-the-age-of-ai/>.
- 352* B. Nimmo, M. Flossman, « Influence et cyberopérations : une mise à jour » (OpenAI, 2024) ; https://cdn.openai.com/threat-intelligence-reports/influence-and-cyber-operations-an-update_October-2024.pdf.
- 353 Agence des projets de recherche avancée de défense, AlxCC (2024) ; <https://aicyberchallenge.com/>.
- 354 H. Ruan, Y. Zhang, A. Roychoudhury, SpecRover : Extraction d'intention de code via des LLM, arXiv [cs.SE] (2024) ; <http://arxiv.org/abs/2408.02232>.
- 355 NT Islam, J. Khoury, A. Seong, E. Bou-Harb, P. Najafirad, Améliorer la sécurité du code source avec les LLM : démystifier les défis et générer des réparations fiables, arXiv [cs.CR] (2024) ; <http://arxiv.org/abs/2409.00571>.
- 356 X. Du, G. Zheng, K. Wang, J. Feng, W. Deng, M. Liu, B. Chen, X. Peng, T. Ma, Y. Lou, Vul-RAG : Améliorer la vulnérabilité basée sur le LLM Détection via Knowledge-Level RAG, arXiv [cs.SE] (2024) ; <http://arxiv.org/abs/2406.11147>.
- 357* M. Allamanis, M. Arjovsky, C. Blundell, L. Buesing, M. Brand, S. Glazunov, D. Maier, P. Maniatis, G. Marinho, H. Michalewski, K. Sen, C. Sutton, V. Tulsyan, M. Vanotti, T. Weber, D. Zheng, De la sieste au grand sommeil : utiliser de grands modèles de langage pour détecter les vulnérabilités dans le code du monde réel (2024) ; <https://googleprojectzero.blogspot.com/2024/10/from-naptime-to-big-sleep.html>.
- 358 AK Zhang, N. Perry, R. Dulepet, J. Ji, J. W. Lin, E. Jones, C. Menders, G. Hussein, S. Liu, D. Jasper, P. Peetathawatchai, A. Glenn, V. Sivashankar, D. Zamoshchin, L. Glikbarg, D. Askaryar, M. Yang, ... P. Liang, Cybench : un cadre d'évaluation des capacités et des risques de cybersécurité des modèles linguistiques, arXiv [cs.CR] (2024) ; <http://arxiv.org/abs/2408.08926>.
- 359 D. Ristea, V. Mavroudis, C. Hicks, Analyse comparative d'OpenAI o1 en cybersécurité, arXiv [cs.CR] (2024) ; <http://arxiv.org/abs/2410.21939>.
- 360 J. Gennari, S.-H. Lau, S. Perl, J. Parish, G. Sastry, « Considérations pour l'évaluation de grands modèles de langage pour « Tâches de cybersécurité » (Université Carnegie Mellon, 2024) ; <https://insights.sei.cmu.edu/library/considerations-for-evaluating-large-language-models-for-cybersecurity-tasks/>.
- 361 M. Shao, B. Chen, S. Jancheska, B. Dolan-Gavitt, S. Garg, R. Karri, M. Shafique, Une évaluation empirique des LLM pour résoudre les défis de sécurité offensive, arXiv [cs.CR] (2024) ; <http://arxiv.org/abs/2402.11814>.
- 362* J. Xu, JW Stokes, G. McDonald, X. Bai, D. Marshall, S. Wang, A. Swaminathan, Z. Li, AutoAttacker : un système guidé par un modèle de langage de grande taille pour mettre en œuvre des cyberattaques automatiques, arXiv [cs.CR] (2024) ; <http://arxiv.org/abs/2403.01038>.
- 363 R. Fang, R. Bindu, A. Gupta, Q. Zhan, D. Kang, Les équipes d'agents LLM peuvent exploiter les vulnérabilités zero-day, arXiv [cs.MA] (2024) ; <http://arxiv.org/abs/2406.01637>.
- 364 T. Abramovich, M. Udeshi, M. Shao, K. Lieret, H. Xi, K. Milner, S. Jancheska, J. Yang, CE Jimenez, F. Khorrami, P. Krishnamurthy, B. Dolan-Gavitt, M. Shafique, K. Narasimhan, R. Karri, O. Press, EnIGMA : agent de modèle génératif interactif amélioré pour les défis CTF, arXiv [cs.AI] (2024) ; <http://arxiv.org/abs/2409.16165>.
- 365 G. Deng, Y. Liu, V. Mayoral-Vilches, P. Liu, Y. Li, Y. Xu, T. Zhang, Y. Liu, M. Pinzger, S. Rass, « PentestGPT : évaluation et exploitation des grandes Modèles de langage pour les tests d'intrusion automatisés » au 33e USENIX Security Symposium (USENIX Security 24) (USENIX Association, Philadelphie, PA, 2024), pages 847 à 864 ; <https://www.usenix.org/conference/usenixsecurity24/presentation/deng>.
- 366* S. Glazunov, M. Brand, Google Project Zero, « Projet Naptime : évaluation des capacités de sécurité offensives des grands modèles de langage » (Google Project Zero, 2024) ; <https://googleprojectzero.blogspot.com/2024/06/project-naptime.html>.
- 367 J. Walden, « L'impact d'un événement de sécurité majeur sur un projet Open Source : le cas d'OpenSSL » dans Actes de la 17e Conférence internationale sur les référentiels de logiciels miniers (ACM, New York, NY, États-Unis,

- 2020); <https://doi.org/10.1145/3379597.3387465>.
- 368 G. Kokolakis, A. Moschos, AD Keromytis, « Exploiter la puissance des LLM à usage général dans la conception de chevaux de Troie matériels » dans *Lecture Notes in Computer Science* (Springer Nature Switzerland, Cham, 2024) Notes de cours en informatique, pp. 176–194 ; https://doi.org/10.1007/978-3-031-61486-6_11.
- 369 JP Farwell, R. Rohozinski, Stuxnet et l'avenir de la cyberguerre. *Survival* 53, 23–40 (2011) ; <https://doi.org/10.1080/00396338.2011.555586>.
- 370 D. Saha, S. Tarek, K. Yahyaei, SK Saha, J. Zhou, M. Tehranipoor, F. Farahmandi, LLM pour la sécurité des systèmes sur puce : un changement de paradigme. *IEEE Access* 12, 155498–155521 (2024) ; <https://doi.org/10.1109/ACCESS.2024.3427369>.
- 371* Amazon, Qu'est-ce qu'AWS CloudTrail ? (2024) ; <https://docs.aws.amazon.com/awsccloudtrail/latest/userguide/cloudtrail-user-guide.html>.
- 372* P. Kanuparth, A. Dalakoti, S. Kamath, Débogage de l'IA chez Meta avec HawkEye, *Engineering chez Meta* (2023) ; <https://engineering.fb.com/2023/12/19/data-infrastructure/hawkeye-ai-debugging-meta/>.
- 373 MC Horowitz, P. Scharre, A. Velez-Green, Un avenir nucléaire stable ? L'impact des systèmes autonomes et Intelligence artificielle, *arXiv [cs.CY]* (2019) ; <http://arxiv.org/abs/1912.05291>.
- 374 AE Chu, T. Lu, P.-S. Huang, Sparks of Function par de Novo Protein Design. *Nature Biotechnology* 42, 203–215 (2024) ; <https://doi.org/10.1038/s41587-024-02133-2>.
- 375 Robert F. Service, Les outils d'IA déclenchent une explosion de protéines de conception. *Science* 386, 260–261 (2024) ; <https://doi.org/10.1126/science.adt9024>.
- 376 C. Li, G. Ye, Y. Jiang, Z. Wang, H. Yu, M. Yang, L'intelligence artificielle dans la lutte contre les maladies infectieuses : un rôle transformateur. *Journal of Medical Virology* 96, e29355 (2024) ; <https://doi.org/10.1002/jmv.29355>.
- 377 Académie royale des sciences de Suède, Prix Nobel de chimie 2024. (2024) ; <https://www.nobelprize.org/uploads/2024/10/press-chemistryprize2024-3.pdf>.
- 378 V. Pitschmann, Z. Hon, Les drogues comme armes chimiques : passé et perspectives. *Toxics* 11, 52 (2023) ; <https://doi.org/10.3390/toxics11010052>.
- 379 National Research Council, « Biosécurité et recherche à double usage dans les sciences de la vie » dans *Science et sécurité dans un monde post-11 septembre : un rapport basé sur des discussions régionales entre les communautés scientifiques et de sécurité* (National Academies Press, Washington, DC, DC, 2007) ; <https://doi.org/10.17226/1213>.
- 380 S. Ben Ouagrham-Gormley, *Obstacles aux armes biologiques : les défis de l'expertise et de l'organisation pour le développement des armes* (Cornell University Press, 2014) ; <https://www.cornellpress.cornell.edu/book/9780801452888/barriers-to-bioweapons/#bookTabs=1>.
- 381 J. Reville, C. Jefferson, Tacit Knowledge and the Biological Weapons Regime. *Science & Public Policy* 41, 597–610 (2014) ; <https://doi.org/10.1093/scipol/sct090>.
- 382 SR Carter, N. Wheeler, S. Chwalek, C. Isaac, JM Yassif, « La convergence de l'intelligence artificielle et des sciences de la vie : protéger la technologie, repenser la gouvernance et prévenir les catastrophes » (Nuclear Threat Initiative, 2023) ; https://www.nti.org/wp-content/uploads/2023/10/NTIBIO_AI_FINAL.pdf.
- 383 J. Smith, S. Rose, R. Moulange, C. Nelson, « Comment le gouvernement britannique devrait-il faire face au risque d'utilisation abusive des outils biologiques basés sur l'IA » (Centre pour la résilience à long terme, 2024) ; <https://www.longtermresilience.org/wp-content/uploads/2024/07/How-the-UK-Government-should-address-the-misuse-risk-from-AI-Enabled-biological-tools-BTs-Website-Copy.pdf>.
- 384 B. Drexel, C. Withers, « L'IA et l'évolution des risques biologiques pour la sécurité nationale : capacités, seuils et interventions » (CNAS, 2024) ; <https://www.cnas.org/publications/reports/ai-and-the-evolution-of-biological-national-security-risks>.
- 385 M. Dybul, « La biosécurité à l'ère de l'IA : déclaration du président » (Helena, 2024) ; <https://www.helenabiosecurity.org/>.
- 386* T. Hayes, R. Rao, H. Akin, NJ Sofroniew, D. Oktay, Z. Lin, R. Verkuil, VQ Tran, J. Deaton, M. Wiggert, R. Badkundri, I. Shafkat, J. Gong, A. Derry, RS Molina, N. Thomas, Y. Khan, ... A. Rives, Simulation de 500 millions d'années d'évolution avec un modèle de langage, *bioRxiv [préimpression]* (2024) ; <https://doi.org/10.1101/2024.07.01.600583>.
- 387* V. Zambaldi, D. La, AE Chu, H. Patani, AE Danson, TOC Kwan, T. Frerix, RG Schneider, D. Saxton, A. Thillaisundaram, Z. Wu, I. Moraes, O. Lange, E. Papa, G. Stanton, V. Martin, S. Singh, ... J. Wang, « Conception de novo de liants protéiques à haute affinité avec AlphaProteo » (Google DeepMind, 2024) ; <https://deepmind.google/discover/blog/alphaproteo-generates-novel-proteins-for-biology-and-health-research/>.
- 388 Frontier Model Forum, Mise à jour des progrès : faire progresser la sécurité de l'IA de pointe en 2024 et au-delà, *Frontier Model Forum* (2024) ; <https://www.frontiermodelforum.org/updates/progress-update-advancing-frontier-ai-safety-in->

2024 et au-delà/.

- 389 AlxBio Global Forum, « Livre blanc : Structure et objectifs du Forum mondial AlxBio » (NTI, 2024) ; https://www.nti.org/wp-content/uploads/2024/07/Al_Bio-Global-Forum-Structure-et-Objectifs_Livre-Blanc.pdf .
- 390 NN Thadani, S. Gurev, P. Notin, N. Youssef, NJ Rollins, D. Ritter, C. Sander, Y. Gal, DS Marks, Apprendre de Données pré-pandémiques pour prévoir l'échappement viral. *Nature* 622, 818–825 (2023) ; <https://doi.org/10.1038/s41586-023-06617-0>.
- 391 EH Soice, R. Rocha, K. Cordova, M. Specter, KM Esvelt, Les grands modèles de langage peuvent-ils démocratiser l'accès à la biotechnologie à double usage ?, *arXiv [cs.CY]* (2023) ; <http://arxiv.org/abs/2306.03809>.
- 392 N. Li, A. Pan, A. Gopal, S. Yue, D. Berrios, A. Gatti, JD Li, A.-K. Dombrowski, S. Goel, L. Phan, G. Mukobi, N. Helm-Burger, R. Lababidi, L. Justen, AB Liu, M. Chen, I. Barrass, ... D. Hendrycks, The WMDP Benchmark : Mesurer et réduire les utilisations malveillantes grâce au désapprentissage, *arXiv [cs.LG]* (2024) ; <http://arxiv.org/abs/2403.03218>.
- 393 CA Mouton, C. Lucas, E. Guest, « Les risques opérationnels de l'IA dans les attaques biologiques à grande échelle : résultats d'une étude Red-Team » (RAND Corporation, 2024) ; https://www.rand.org/pubs/research_reports/RRA2977-2.html.
- 394* T. Patwardhan, K. Liu, T. Markov, N. Chowdhury, D. Leet, N. Cone, C. Maltbie, J. Huizinga, C. Wainwright, S. (froggi) Jackson, S. Adler, R. Casagrande, A. Madry, « Construire un système d'alerte précoce pour la création de menaces biologiques assistée par LLM » (OpenAI, 2024) ; <https://openai.com/research/building-an-early-warning-system-for-llm-aided-biological-threat-creation> .
- 395 B.J. Wittmann, T. Alexanian, C. Bartling, J. Beal, A. Clore, J. Diggans, K. Flyngolts, B.T. Gemler, T. Mitchell, ST Murphy, NE Wheeler, E. Horvitz, Vers un criblage résilient à l'IA des commandes de synthèse d'acides nucléiques : processus, résultats et recommandations, *bioRxiv [préimpression]* (2024) ; <https://doi.org/10.1101/2024.12.02.626439>.
- 396 NR Bennett, B. Coventry, I. Goreshnik, B. Huang, A. Allen, D. Vafeados, YP Peng, J. Dauparas, M. Baek, L. Stewart, F. DiMaio, S. De Munck, SN Savvides, D. Baker, Amélioration de la conception de liants protéiques de Novo grâce au Deep Learning. *Nature Communications* 14, 2625 (2023); <https://doi.org/10.1038/s41467-023-38328-5>.
- 397 M. Crowley, L. Shang, M. Dando, Préserver la norme contre les armes chimiques : une initiative de la société civile pour la 4e Conférence d'examen de la Convention sur les armes chimiques de 2018. *Futures* 102, 125–133 (2018) ; <https://doi.org/10.1016/j.futures.2018.01.006>.
- 398 F. Urbina, F. Lentzos, C. Invernizzi, S. Ekins, Double utilisation de la découverte de médicaments alimentée par l'intelligence artificielle. *Nature Machine Intelligence* 4, 189–191 (2022) ; <https://doi.org/10.1038/s42256-022-00465-9>.
- 399 M. Guo, Z. Li, X. Deng, D. Luo, J. Yang, Y. Chen, W. Xue, ConoDL : Un cadre d'apprentissage profond pour une génération rapide et prédiction des conotoxines, *bioRxiv [préimpression]* (2024) ; <https://doi.org/10.1101/2024.09.27.614001>.
- 400* 310.ai, GenAI + BIO : La nature n'a pas eu le temps, nous avons des GPU (2024) ; <https://310.ai/>.
- 401* Asimov, Kernel : Logiciel de CAO pour l'ingénierie biologique (2024) ; <https://www.asimov.com/kernel>.
- 402 A. M Bran, S. Cox, O. Schilter, C. Baldassari, AD White, P. Schwaller, Augmentation des grands modèles de langage avec Outils de chimie. *Nature Machine Intelligence* 6, 525–535 (2024) ; <https://doi.org/10.1038/s42256-024-00832-8>.
- 403 J. Goldblat, La Convention sur les armes biologiques : un aperçu. *Revue internationale de la Croix-Rouge* 37, 251–265 (1997); <https://doi.org/10.1017/s0020860400084679>.
- 404 G. Gonzalez-Isunza, MZ Jawaid, P. Liu, DL Cox, M. Vazquez, J. Arsuaga, Utilisation de l'apprentissage automatique pour détecter les coronavirus potentiellement infectieux pour les humains. *Scientific Reports* 13, 9319 (2023) ; <https://doi.org/10.1038/s41598-023-35861-7>.
- 405 M. Wardeh, MSC Blagrove, KJ Sharkey, M. Baylis, Diviser pour mieux régner : l'apprentissage automatique intègre les traits des mammifères et des virus aux caractéristiques du réseau pour prédire les associations virus-mammifères. *Nature Communications* 12, 3954 (2021) ; <https://doi.org/10.1038/s41467-021-24085-w>.
- 406 S. Rose, R. Moulange, J. Smith, C. Nelson, « L'impact à court terme de l'IA sur l'utilisation abusive de la biologie » (Centre pour la résilience à long terme, 2024) ; <https://www.longtermresilience.org/reports/the-near-term-impact-of-ai-on-biological-misuse/>.
- 407 J. Frazer, P. Notin, M. Dias, A. Gomez, JK Min, K. Brock, Y. Gal, DS Marks, Prédiction des variantes de maladies avec Deep Modèles génératifs de données évolutives. *Nature* 599, 91–95 (2021) ; <https://doi.org/10.1038/s41586-021-04043-8>.
- 408 JB Sandbrink, EC Alley, MC Watson, GD Koblenz, KM Esvelt, Insidious Insights : implications de l'ingénierie des vecteurs viraux pour l'amélioration des agents pathogènes. *Thérapie génique* 30, 407–410 (2023) ; <https://doi.org/10.1038/s41434-021-00312-3>.
- 409 J. Kaiser, Exclusif : Des expériences controversées qui pourraient rendre la grippe aviaire plus risquée sont sur le point de reprendre, selon American

- Association pour l'avancement de la science (2021) ; <https://www.science.org/content/article/exclusive-controversial-experiments-make-bird-flu-more-risky-poised-resume> .
- 410 J. Pannu, D. Bloomfield, A. Zhu, R. MacKnight, G. Gomes, A. Cicero, T. Inglesby, Priorisation des capacités biologiques à conséquences élevées dans les évaluations des modèles d'intelligence artificielle, arXiv [cs.CY] (2024) ; <http://dx.doi.org/10.2139/ssrn.4873106>.
- 411 E. Appleton, C. Madsen, N. Roehner, D. Densmore, Automatisation de la conception en biologie synthétique. *Cold Spring Harbor Perspectives in Biology* 9 (2017) ; <https://doi.org/10.1101/cshperspect.a023978>.
- 412 Organisation de coopération et de développement économiques, Intelligence artificielle dans la science : défis, Opportunités et avenir de la recherche (OCDE, Paris, 2023) ; https://www.oecd-ilibrary.org/science-and-technology/artificial-intelligence-in-science_a8d820bd-fr .
- 413 C. Nelson, S. Rose, « Comprendre le développement d'armes biologiques facilité par l'IA » (Centre de recherche sur les armes biologiques à long terme) Résilience, 2023) ; <https://www.longtermresilience.org/reports/understanding-risks-at-the-intersection-of-ai-and-bio/> .
- 414 Z. Wu, SBJ Kan, RD Lewis, BJ Wittmann, FH Arnold, Évolution dirigée des protéines assistée par apprentissage automatique avec bibliothèques combinatoires. *Actes de l'Académie nationale des sciences des États-Unis d'Amérique* 116, 8852–8858 (2019) ; <https://doi.org/10.1073/pnas.1901979116>.
- 415 DA Boiko, R. MacKnight, B. Kline, G. Gomes, Recherche chimique autonome avec de grands modèles de langage. *Nature* 624, 570–578 (2023) ; <https://doi.org/10.1038/s41586-023-06792-0>.
- 416 A. Stephenson, L. Lastra, B. Nguyen, Y.-J. Chen, J. Nivala, L. Ceze, K. Strauss, Automatisation physique du laboratoire en biologie synthétique. *ACS Synthetic Biology* 12, 3156–3169 (2023) ; <https://doi.org/10.1021/acssynbio.3c00345>.
- 417 JT Rapp, BJ Bremer, PA Romero, Laboratoires autonomes pour naviguer de manière autonome dans le paysage de la forme physique des protéines. *Nature Chemical Engineering* 1, 97–107 (2024) ; <https://doi.org/10.1038/s44286-023-00002-4>.
- 418 A. Casas, M. Bultelle, R. Kitney, Une approche d'ingénierie biologique pour le flux de travail automatisé et la bioconception. *Biologie synthétique* 9, ysae009 (2024) ; <https://doi.org/10.1093/synbio/ysae009>.
- 419 D. Sun, W. Gao, H. Hu, S. Zhou, Pourquoi 90 % du développement clinique de médicaments échoue et comment l'améliorer ? *Acta Pharmaceutica Sinica* B 12, 3049-3062 (2022) ; <https://doi.org/10.1016/j.apsb.2022.02.002>.
- 420 Forum sur les neurosciences et les troubles du système nerveux, Conseil sur la politique des sciences de la santé, Institut de médecine, « Défis du développement de médicaments » dans *Améliorer et accélérer le développement thérapeutique pour les troubles du système nerveux : résumé de l'atelier* (National Academies Press (États-Unis), 2014) ; <https://www.ncbi.nlm.nih.gov/books/NBK195047/>.
- 421 KH Sumida, R. Núñez-Franco, I. Kalvet, SJ Pellock, BIM Wicky, LF Milles, J. Dauparas, J. Wang, Y. Kipnis, N. Jameson, A. Kang, J. De La Cruz, B. Sankaran, AK Bera, G. Jiménez-Osés, D. Baker, Amélioration de l'expression, de la stabilité et de la fonction des protéines avec ProteinMPNN. *Journal of the American Chemical Society* 146, 2054–2061 (2024) ; <https://doi.org/10.1021/jacs.3c10941>.
- 422 M. Wehrs, D. Tanjore, T. Eng, J. Lievense, TR Pray, A. Mukhopadhyay, Ingénierie de microbes de production robustes pour la culture à grande échelle. *Tendances en microbiologie* 27, 524–537 (2019) ; <https://doi.org/10.1016/j.tim.2019.01.006>.
- 423 J. Jiang, H.-H. Peng, Z. Yang, X. Ma, S. Sahakijpipjam, C. Moon, D. Ouyang, RO Williams Iii, Les applications de l'apprentissage automatique (ML) dans la conception de poudre sèche pour inhalation en utilisant la technologie de congélation en couche mince. *Journal international de pharmacie* 626, 122179 (2022) ; <https://doi.org/10.1016/j.ijpharm.2022.122179>.
- 424 TR Sosnowski, Vers un ciblage plus précis des aérosols inhalés sur différentes zones du système respiratoire. *Pharmaceutics* 16, 97 (2024) ; <https://doi.org/10.3390/pharmaceutics16010097>.
- 425 Département des sciences, de l'innovation et de la technologie, Institut de sécurité de l'IA, « Évaluations avancées de l'IA à l'AIISI : mai « Mise à jour » (GOV.UK, 2024) ; <https://www.aisi.gov.uk/work/advanced-ai-evaluations-may-update>.
- 426* Anthropic, Réflexions sur notre politique de mise à l'échelle responsable (2024) ; <https://www.anthropic.com/news/reflections-sur-notre-politique-de-mise-à-l'échelle-responsable>.
- 427 G. Lewis, P. Millett, A. Sandberg, A. Snyder-Beattie, G. Gronvall, Information Hazards in Biotechnology. *Analyse des risques : une publication officielle de la Society for Risk Analysis* 39, 975–981 (2019) ; <https://doi.org/10.1111/risa.13235>.
- 428 SR Carter, S. Curtis, C. Emerson, J. Gray, IC Haydon, A. Hebbeler, C. Qureshi, N. Randolph, A. Rives, AL Stuart, Responsible AI X Biodesign : valeurs communautaires, principes directeurs et engagements pour le développement responsable de l'IA pour la conception de protéines (2024) ; <https://responsiblebiodesign.ai/>.
- 429 NTI | bio, « Projet d'agenda de recherche pour la sauvegarde des capacités IA-Bio » (NTI, 2024) ; <https://www.nti.org/wp-content/uploads/2024/06/Research-Agenda-for-Safeguarding-AI-Bio-Capabilities.pdf> .
- 430 E. Nguyen, M. Poli, MG Durrant, AW Thomas, B. Kang, J. Sullivan, MY Ng, A. Lewis, A. Patel, A. Lou, S. Ermon, S.

- A. Baccus, T. Hernandez-Boussard, C. Re, PD Hsu, BL Hie, Modélisation et conception de séquences de l'échelle moléculaire à l'échelle du génome avec Evo, bioRxiv [préimpression] (2024) ; <https://doi.org/10.1101/2024.02.27.582234>.
- 431 J. Cheng, G. Novati, J. Pan, C. Bycroft, A. Žemgulytė, T. Applebaum, A. Pritzel, LH Wong, M. Zielinski, T. Sargeant, R. G. Schneider, AW Senior, J. Jumper, D. Hassabis, P. Kohli, Ž. Avsec, Prédiction précise de l'effet des variants faux-sens à l'échelle du protéome avec AlphaMissense. *Science (New York, NY)* 381, eadg7492 (2023) ; <https://doi.org/10.1126/science.adg7492>.
- 432 SR Carter, NE Wheeler, C. Isaac, JM Yassif, « Développement de garde-fous pour les outils de bioconception de l'IA » (Nuclear Threat Initiative, 2024) ; <https://www.nti.org/analysis/articles/developing-guardrails-for-ai-biodesign-tools/>.
- 433 SA Dip, UA Shuvo, T. Chau, H. Song, P. Choi, X. Wang, L. Zhang, PathoLM : Identification de la pathogénicité à partir de la séquence d'ADN via le modèle de fondation du génome, arXiv [cs.CL] (2024) ; <http://arxiv.org/abs/2406.13133>.
- 434 K. Workman, Ingénierie des AAV avec Evo et AlphaFold, LatchBio (2024) ; <https://blog.latch.bio/p/engineering-aavs-with-evo-and-alpha-fold>.
- 435 D. Bloomfield, J. Pannu, AW Zhu, MY Ng, A. Lewis, E. Bendavid, SM Asch, T. Hernandez-Boussard, A. Cicero, T. Inglesby, IA et biosécurité : la nécessité d'une gouvernance. *Science (New York, NY)* 385, 831–833 (2024) ; <https://doi.org/10.1126/science.adq1977>.
- 436 Y. Zhang, M. Yasunaga, Z. Zhou, JZ HaoChen, J. Zou, P. Liang, S. Yeung, « Au-delà de la mise à l'échelle positive : comment la négation affecte les tendances de mise à l'échelle des modèles linguistiques » dans *Conclusions de l'Association pour la linguistique computationnelle : ACL 2023*, A. Rogers, J. Boyd-Graber, N. Okazaki, éd. (Association pour la linguistique computationnelle, 2023), pp. 7479–7498 ; <https://doi.org/10.18653/v1/2023.findings-acl.472>.
- 437 A. Mallen, A. Asai, V. Zhong, R. Das, D. Khashabi, H. Hajishirzi, « Quand ne pas faire confiance aux modèles linguistiques : étude de l'efficacité des mémoires paramétriques et non paramétriques » dans *Actes de la 61e réunion annuelle de l'Association for Computational Linguistics (Volume 1 : Longs articles)*, A. Rogers, J. Boyd-Graber, N. Okazaki, éd. (Association de linguistique computationnelle, Toronto, Canada, 2023), pp. 9802–9822 ; <https://doi.org/10.18653/v1/2023.acl-long.546>.
- 438 S. Santurkar, E. Durmus, F. Ladhak, C. Lee, P. Liang, T. Hashimoto, « À qui les modèles linguistiques reflètent-ils les opinions ? » dans *Actes de la 40e Conférence internationale sur l'apprentissage automatique (JMLR, Honolulu, Hawaï, États-Unis, 2023)* vol. 202 de ICML'23, pp. 29971–30004 ; <https://proceedings.mlr.press/v202/santurkar23a.html>.
- 439 L. Weidinger, J. Uesato, M. Rauh, C. Griffin, P.-S. Huang, J. Mellor, A. Glaese, M. Cheng, B. Balle, A. Kasirzadeh, C. Biles, S. Brown, Z. Kenton, W. Hawkins, T. Stepleton, A. Birhane, LA Hendricks, ... I. Gabriel, « Taxonomie des risques posés par les modèles linguistiques » dans *Actes de la conférence 2022 de l'ACM sur l'équité, la responsabilité et la transparence (FAccT '22)* (Association for Computing Machinery, New York, NY, États-Unis, 2022), pp. 214–229 ; <https://doi.org/10.1145/3531146.3533088>.
- 440* M. Chen, J. Twarek, H. Jun, Q. Yuan, HP de Oliveira Pinto, J. Kaplan, H. Edwards, Y. Burda, N. Joseph, G. Brockman, A. Ray, R. Puri, G. Krueger, M. Petrov, H. Khlaaf, G. Sastry, P. Mishkin, ... W. Zaremba, Évaluation de grands modèles de langage formés au code, arXiv [cs.LG] (2021) ; <http://arxiv.org/abs/2107.03374>.
- 441 S. Nguyen, HM Babe, Y. Zi, A. Guha, CJ Anderson, MQ Feldman, « Comment les programmeurs débutants et les LLM en code se (mal)interprètent » dans *Actes de la conférence CHI sur les facteurs humains dans les systèmes informatiques (CHI '24)* (Association for Computing Machinery, New York, NY, États-Unis, 2024), pp. 1–26 ; <https://doi.org/10.1145/3613904.3642706>.
- 442 F. Cassano, L. Li, A. Sethi, N. Shinn, A. Brennan-Jones, J. Ginesin, E. Berman, G. Chakhnashvili, A. Lozhkov, CJ Anderson, A. Guha, Can It Edit ? Évaluation de la capacité des grands modèles de langage à suivre les instructions d'édition de code, arXiv [cs.SE] (2023) ; <http://arxiv.org/abs/2312.12450>.
- 443 R. Pan, AR Ibrahimzada, R. Krishna, D. Sankar, LP Wassi, M. Merler, B. Sobolev, R. Pavuluri, S. Sinha, R. Jabbarvand, « Lost in Translation : une étude des bugs introduits par les grands modèles de langage lors de la traduction de code » dans *les actes de la 46e conférence internationale IEEE/ACM sur l'ingénierie logicielle (ICSE '24)* (Association for Computing Machinery, New York, NY, États-Unis, 2024), pp. 1–13 ; <https://doi.org/10.1145/3597503.3639226>.
- 444 N. Perry, M. Srivastava, D. Kumar, D. Boneh, « Les utilisateurs écrivent-ils du code moins sécurisé avec des assistants IA ? » dans *Actes de la conférence ACM SIGSAC 2023 sur la sécurité informatique et des communications (ACM, New York, NY, États-Unis, 2023)*, pp. 2785–2799 ; <https://doi.org/10.1145/3576915.3623157>.
- 445 A. Perlman, Les implications de ChatGPT pour les services juridiques et la société, *The Practice* (2023) ; <https://clp.law.harvard.edu/knowledge-hub/magazine/issues/generative-ai-in-the-legal-profession/the-implications-of-chatgpt-for-legal-services-and-society/>.
- 446 E. Martinez, Réévaluation des performances du GPT-4 à l'examen du barreau. *Intelligence artificielle et droit* (2024) ; <https://doi.org/10.1007/s10506-024-09396-9>.
- 447 Eastern District of Texas, Tribunal de district des États-Unis, Mémoire et ordonnance dans l'affaire 1 : 23-Cv-00281-MAC. (2024) ;

- <https://www.courthousenews.com/wp-content/uploads/2024/11/un-avocat-sanctionne-pour-avoir-utilise-des-hallucinations-d-intelligence-intellectuelle.pdf> .
- 448 JA Omiye, JC Lester, S. Spichak, V. Rotemberg, R. Daneshjou, Les grands modèles linguistiques propagent la médecine fondée sur la race. *Npj Digital Medicine* 6, 1–4 (2023) ; <https://doi.org/10.1038/s41746-023-00939-z>.
- 449 TH Kung, M. Cheatham, A. Medenilla, C. Sillos, L. De Leon, C. Elepaño, M. Madriaga, R. Aggabao, G. Diaz-Candido, J. Maningo, V. Tseng, Performance de ChatGPT sur USMLE : potentiel pour l'enseignement médical assisté par l'IA à l'aide de grands modèles linguistiques. *PLOS Digital Health* 2, e0000198 (2023) ; <https://doi.org/10.1371/journal.pdig.0000198>.
- 450 K. Singhal, S. Azizi, T. Tu, SS Mahdavi, J. Wei, HW Chung, N. Scales, A. Tanwani, H. Cole-Lewis, S. Pfohl, P. Payne, M. Seneviratne, P. Gamble, C. Kelly, A. Babiker, N. Schärli, A. Chowdhery, ... V. Natarajan, Les grands modèles linguistiques codent les connaissances cliniques. *Nature* 620, 172–180 (2023) ; <https://doi.org/10.1038/s41586-023-06291-2>.
- 451 J. Tan, H. Westermann, K. Benyekhlef, « ChatGPT en tant qu'avocat artificiel ? » dans *Atelier sur l'intelligence artificielle pour l'accès à la justice (AI4AJ 2023)* (Actes de l'atelier CEUR, Braga, Portugal, 2023) ; <https://ceur-ws.org/Vol-3435/short2.pdf>.
- 452 JLM Brand, Le chatbot d'Air Canada illustre les problèmes persistants d'agence et d'écart de responsabilité pour l'IA. *AI & Society*, 1–3 (2024) ; <https://doi.org/10.1007/s00146-024-02096-7>.
- 453* Z. Yuan, H. Yuan, C. Tan, W. Wang, S. Huang, Dans quelle mesure les grands modèles de langage sont-ils performants dans les tâches arithmétiques ? *arXiv [cs.CL]* (2023) ; <http://arxiv.org/abs/2304.02015>.
- 454 Z. Wang, « CausalBench : une référence complète pour évaluer les capacités de raisonnement causal des grandes entreprises Modèles linguistiques » dans les actes du 10e atelier SIGHAN sur le traitement de la langue chinoise (SIGHAN-10) (2024), pp. 143–151 ; <https://aclanthology.org/2024.sighan-1.17.pdf>.
- 455 X. Yin, J. Jiang, L. Yang, X. Wan, History Matters : Édition des connaissances temporelles dans un grand modèle linguistique. Actes de la ... Conférence AAAI sur l'intelligence artificielle. Conférence AAAI sur l'intelligence artificielle 38, 19413–19421 (2024) ; <https://doi.org/10.1609/aaai.v38i17.29912>.
- 456 ID Raji, IE Kumar, A. Horowitz, A. Selbst, « L'illusion de la fonctionnalité de l'IA » dans les actes de l'ACM 2022 Conférence sur l'équité, la responsabilité et la transparence (FAccT '22) (Association for Computing Machinery, New York, NY, États-Unis, 2022), pp. 959–972 ; <https://doi.org/10.1145/3531146.3533158>.
- 457 B. Vidgen, A. Agrawal, AM Ahmed, V. Akinwande, N. Al-Nuaimi, N. Alfaraj, E. Alhajjar, L. Aroyo, T. Bavalatti, M. Bartolo, B. Bili-Hamelin, K. Bollacker, R. Bomassani, MF Boston, S. Campos, K. Chakra, C. Chen, ... J. Vanschoren, Présentation de la version 0.5 du benchmark de sécurité de l'IA de MLCommons, *arXiv [cs.CL]* (2024) ; <http://arxiv.org/abs/2404.12241>.
- 458 P. Guldemann, A. Spiridonov, R. Staab, N. Jovanović, M. Vero, V. Vechev, A. Gueorguieva, M. Balunović, N. Konstantinov, P. Bielik, P. Tsankov, M. Vechev, Cadre COMPL-AI : une suite d'interprétation technique et d'analyse comparative LLM pour la loi européenne sur l'intelligence artificielle, *arXiv [cs.CL]* (2024) ; <http://arxiv.org/abs/2410.07959>.
- 459 Observatoire des politiques de l'OCDE en matière d'IA, *OECD AI Incidents Monitor (AIM)* (2024) ; <https://oecd.ai/en/incidents>.
- 460 A. Wei, N. Haghtalab, J. Steinhardt, « Jailbreaké : comment la formation à la sécurité LLM échoue-t-elle ? » dans 37e Conférence sur les systèmes de traitement de l'information neuronale (NeurIPS 2023) (La Nouvelle-Orléans, LA, États-Unis, 2023) ; <https://openreview.net/forum?id=jA235JGM09>.
- 461 SMTI Tonmoy, SMM Zaman, V. Jain, A. Rani, V. Rawte, A. Chadha, A. Das, Une étude complète des techniques d'atténuation des hallucinations dans les grands modèles de langage, *arXiv [cs.CL]* (2024) ; <http://arxiv.org/abs/2401.01313>.
- 462 ETH Zurich, INSAIT, LatticeFlow AI, COMPL-AI (2024) ; <https://compl-ai.org/>.
- 463 N. Guha, J. Nyarko, DE Ho, C. Ré, A. Chilton, A. Narayana, A. Chohlas-Wood, A. Peters, B. Waldon, DN Rockmore, D. Zambrano, D. Talisman, E. Hoque, F. Surani, F. Fagan, G. Sarfaty, GM Dickinson, ... Z. Li, « LEGALBENCH : une référence élaborée en collaboration pour mesurer le raisonnement juridique dans les grands modèles linguistiques » dans 37e Conférence internationale sur les systèmes de traitement de l'information neuronale (NeurIPS 2023) (Curran Associates Inc., Red Hook, NY, États-Unis, 2024), pp. 44123–44279 ; <https://doi.org/10.5555/3666122.3668037>.
- 464 R. Xu, Z. Wang, R.-Z. Fan, P. Liu, Analyse comparative des fuites de référence dans les grands modèles de langage, *arXiv [cs.CL]* (2024) ; <http://arxiv.org/abs/2404.18824>.
- 465 S. Longpre, S. Biderman, A. Albalak, H. Schoelkopf, D. McDuff, S. Kapoor, K. Klyman, K. Lo, G. Ilharco, N. San, M. Rauh, A. Skowron, B. Vidgen, L. Weidinger, A. Narayanan, V. Sanh, D. Adelmani, ... L. Soldaini, The Responsible Foundation Model Development Cheatsheet : Un aperçu des outils et des ressources. *Transactions on Machine Learning Research* (2024) ; <https://openreview.net/pdf?id=tH1dQH20eZ>.
- 466 V. Ojewale, R. Steed, B. Vecchione, A. Birhane, ID Raji, Vers une infrastructure de responsabilisation de l'IA : lacunes et opportunités dans les outils d'audit de l'IA, *arXiv [cs.CY]* (2024) ; <http://arxiv.org/abs/2402.17861>.

- 467 N. Guha, CM Lawrence, LA Gailmard, KT Rodolfa, F. Surani, R. Bommasani, ID Raji, M.-F. Cuéllar, C. Honigsberg, P. Liang, DE Ho, La réglementation de l'IA a son propre problème d'alignement : la faisabilité technique et institutionnelle de la divulgation, de l'enregistrement, de l'octroi de licences et de l'audit. *The George Washington Law Review* 92 (2024) ; https://dho.stanford.edu/wp-content/uploads/AI_Regulation.pdf.
- 468 A. Narayanan, S. Kapoor, AI Snake Oil : ce que l'intelligence artificielle peut faire, ce qu'elle ne peut pas faire et comment faire la différence (Princeton University Press, 2024) ; <https://doi.org/10.1515/9780691249643>.
- 469 J. Buolamwini, T. Gebru, « Gender Shades : disparités de précision intersectionnelles dans les genres commerciaux Classification » dans les actes de la 1ère Conférence sur l'équité, la responsabilité et la transparence (FAT/MM '19) (PMLR, 2018), pp. 77–91 ; <https://proceedings.mlr.press/v81/buolamwini18a.html>.
- 470 J. Angwin, J. Larson, L. Kirchner, S. Mattu, Machine Bias, *ProPublica* (2016) ; <https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing>.
- 471 J. Dressel, H. Farid, L'exactitude, l'équité et les limites de la prédiction de la récidive. *Science Advances* 4, eaao5580 (2018) ; <https://doi.org/10.1126/sciadv.aao5580>.
- 472 Z. Obermeyer, B. Powers, C. Vogeli, S. Mullainathan, Analyse des préjugés raciaux dans un algorithme utilisé pour gérer la santé des populations. *Science* 366, 447–453 (2019) ; <https://doi.org/10.1126/science.aax2342>.
- 473 T. Zack, E. Lehman, M. Suzgun, JA Rodriguez, LA Celi, J. Gichoya, D. Jurafsky, P. Szolovits, DW Bates, R.-EE Abdulnour, AJ Butte, E. Alsentzer, Évaluation du potentiel du GPT-4 à perpétuer les préjugés raciaux et sexistes dans les soins de santé : une étude d'évaluation modèle. *The Lancet. Digital Health* 6, e12–e22 (2024) ; [https://doi.org/10.1016/S2589-7500\(23\)00225-X](https://doi.org/10.1016/S2589-7500(23)00225-X).
- 474 F. Bianchi, P. Kalluri, E. Durmus, F. Ladhak, M. Cheng, D. Nozza, T. Hashimoto, D. Jurafsky, J. Zou, A. Caliskan, « La génération de texte en image facilement accessible amplifie les stéréotypes démographiques à grande échelle » dans Actes de la conférence 2023 de l'ACM sur l'équité, la responsabilité et la transparence (FAccT '23) (Association for Computing Machinery, New York, NY, États-Unis, 2023), pp. 1493–1504 ; <https://doi.org/10.1145/3593013.3594095>.
- 475 S. Ghosh, A. Caliskan, « 'personne' == homme occidental à la peau claire et sexualisation des femmes de couleur : « Stéréotypes en diffusion stable » dans les conclusions de l'Association for Computational Linguistics : EMNLP 2023 (Association for Computational Linguistics, Stroudsburg, PA, États-Unis, 2023), pp. 6971–6985 ; <https://doi.org/10.18653/v1/2023.findings-emnlp.465>.
- 476 M. Cheong, E. Abedin, M. Ferreira, R. Reimann, S. Chalson, P. Robinson, J. Byrne, L. Ruppner, M. Alfano, C. Klein, Enquête sur les préjugés sexistes et raciaux dans les mini-images DALL-E. *ACM Journal on Responsible Computing* 1, 1–20 (2024) ; <https://doi.org/10.1145/3649883>.
- 477 JS Park, MS Bernstein, RN Brewer, E. Kamar, MR Morris, « Comprendre la représentation et « Représentativité de l'âge dans les ensembles de données d'IA » dans les actes de la conférence 2021 de l'AAAI/ACM sur l'IA, l'éthique et la société (AIES '21) (Association for Computing Machinery, New York, NY, États-Unis, 2021), pp. 834–842 ; <https://doi.org/10.1145/3461702.3462590>.
- 478 R. Kamikubo, L. Wang, C. Marte, A. Mahmood, H. Kacorri, « Représentativité des données dans les ensembles de données d'accessibilité : une méta-analyse » dans Actes de la 24e conférence internationale ACM SIGACCESS sur les ordinateurs et l'accessibilité (ASSETS '22) (Association for Computing Machinery, New York, NY, États-Unis, 2022), pp. 1–15 ; <https://doi.org/10.1145/3517428.3544826>.
- 479* S. Shankar, Y. Halpern, E. Breck, J. Atwood, J. Wilson, D. Sculley, « Pas de classification sans représentation : évaluation des problèmes de géodiversité dans les ensembles de données ouvertes pour le monde en développement » dans 31e Conférence sur les systèmes de traitement de l'information neuronale (NIPS 2017) Atelier sur l'apprentissage automatique pour le monde en développement (Long Beach, CA, États-Unis, 2017) ; <https://arxiv.org/abs/1711.08536>.
- 480 T. de Vries, I. Misra, C. Wang, L. van der Maaten, « La reconnaissance d'objets fonctionne-t-elle pour tout le monde ? » dans les actes des ateliers de la conférence IEEE/CVF sur la vision par ordinateur et la reconnaissance de formes (CVPR) (Long Beach, CA, États-Unis, 2019) ; https://openaccess.thecvf.com/content_CVPRW_2019/papers/cv4gc/de_Vries_Does_Object_Recognition_Work_for_All_CVPRW_2019_paper.pdf.
- 481 S. Longpre, R. Mahari, A. Chen, N. Obeng-Marnu, D. Sileo, W. Brannon, N. Muennighoff, N. Khazam, J. Kabbara, K. Perisetla, X. Wu, E. Shippole, K. Bollacker, T. Wu, L. Villa, S. Pentland, S. Hooker, The Data Provenance Initiative : un audit à grande échelle des licences et de l'attribution des ensembles de données dans l'IA, *arXiv [cs.CL]* (2023) ; <http://arxiv.org/abs/2310.16787>.
- 482 H. Suresh, J. Guttat, « Un cadre pour comprendre les sources de préjudice tout au long du cycle de vie de l'apprentissage automatique » dans Équité et accès dans les algorithmes, les mécanismes et l'optimisation (ACM, New York, NY, États-Unis, 2021) ; <https://doi.org/10.1145/3465416.3483305>.
- 483* L. Weidinger, J. Mellor, M. Rauh, C. Griffin, J. Uesato, P.-S. Huang, M. Cheng, M. Glaese, B. Balle, A. Kasirzadeh, Z. Kenton, S. Brown, W. Hawkins, T. Stepleton, C. Biles, A. Birhane, J. Haas, ... I. Gabriel, « Risques éthiques et sociaux de

- « Les dommages causés par les modèles linguistiques » (Google DeepMind, 2021) ; <http://arxiv.org/abs/2112.04359>.
- 484 J. Nwatu, O. Ignat, R. Mihalcea, « Comblant la fracture numérique : variation des performances selon les niveaux socio-économiques » dans les actes de la conférence 2023 sur les méthodes empiriques en traitement du langage naturel (EMNLP 2023), H. Bouamor, J. Pino, K. Bali, éd. (Association for Computational Linguistics, Singapour, 2023), pp. 10686–10702 ; <https://doi.org/10.18653/v1/2023.emnlp-main.660>.
- 485 A. Pouget, L. Beyer, E. Bugliarello, X. Wang, AP Steiner, X. Zhai, I. Alabdulmohsin, « No Filter : Cultural and Diversité socioéconomique dans les modèles contrastifs de vision et de langage » dans la 38e conférence annuelle sur les systèmes de traitement de l'information neuronale (NeurIPS 2024) (2024) ; <https://openreview.net/pdf?id=UmW9BYj761>.
- 486 S. Nayak, K. Jain, R. Awal, S. Reddy, S. Van Steenkiste, LA Hendricks, K. Stanczak, A. Agrawal, Analyse comparative des modèles de langage visuel pour la compréhension culturelle (Association pour la linguistique computationnelle, 2024) ; <https://aclanthology.org/2024.emnlp-main.329>.
- 487 D. Agarwal, M. Naaman, A. Vashista, AI Suggestions pour homogénéiser l'écriture vers les styles occidentaux et diminuer Nuances culturelles, arXiv [cs.HC] (2024) ; <http://arxiv.org/abs/2409.11360>.
- 488 N. Shahbazi, Y. Lin, A. Asudeh, HV Jagadish, Biais de représentation dans les données : une enquête sur l'identification et Techniques de résolution. ACM Computing Surveys 55, 293:1–293:39 (2023) ; <https://doi.org/10.1145/3588433>.
- 489 SE Whang, Y. Roh, H. Song, J.-G. Lee, Collecte de données et défis de qualité dans l'apprentissage profond : une perspective d'IA centrée sur les données. The VLDB Journal : Very Large Data Bases : Une publication du VLDB Endowment 32, 791–813 (2023) ; <https://doi.org/10.1007/s00778-022-00775-9>.
- 490 AP Gema, JOJ Leang, G. Hong, A. Devoto, ACM Mancino, R. Saxena, X. He, Y. Zhao, X. Du, MRG Madani, C. Barale, R. McHardy, J. Harris, J. Kaddour, E. van Krieken, P. Minervini, Avons-nous fini avec MMLU ?, arXiv [cs.CL] (2024) ; <http://arxiv.org/abs/2406.04127>.
- 491 Y. Wan, G. Pu, J. Sun, A. Garimella, K.-W. Chang, N. Peng, « "Kelly est une personne chaleureuse, Joseph est un modèle" : les préjugés sexistes dans les lettres de référence générées par les LLM » dans Findings of the Association for Computational Linguistics: EMNLP 2023, H. Bouamor, J. Pino, K. Bali, éd. (Association for Computational Linguistics, Singapour, 2023), pp. 3730–3748 ; <https://doi.org/10.18653/v1/2023.findings-emnlp.243>.
- 492 D. van Niekerk, M. Pérez-Ortiz, J. Shawe-Taylor, D. Orlić, I. Drobnyak, J. Kay, N. Siegel, K. Evans, N. Moorosi, T. Eliassi-Rad, LM Tanczer, W. Holmes, MP Deisenroth, I. Straw, M. Fasli, R. Adams, N. Oliver, ... M. Janicky, « Remettre en question les préjugés systématiques : une enquête sur les préjugés à l'encontre des femmes et des filles dans les grands modèles linguistiques » (UNESCO, IRCAL, 2024) ; <https://ircal.org/project/challenging-systematic-prejudices/>.
- 493 M. Vlasceanu, DM Amodio, Propagation de l'inégalité sociale entre les sexes par les algorithmes de recherche sur Internet. Actes de l'Académie nationale des sciences 119, e2204529119 (2022) ; <https://doi.org/10.1073/pnas.2204529119>.
- 494 S. Sterlie, N. Weng, A. Feragen, Généralisation de l'équité aux modèles de langage génératif via la reformulation de non-Critères de discrimination, arXiv [cs.CL] (2024) ; <http://arxiv.org/abs/2403.08564>.
- 495 T. Sandoval-Martin, E. Martínez-Sanzo, Perpétuation des préjugés sexistes dans la représentation visuelle des professions dans les outils d'IA générative DALL-E et Bing Image Creator. Sciences sociales (Bâle, Suisse) 13, 250 (2024) ; <https://doi.org/10.3390/socsci13050250>.
- 496 L. Sun, M. Wei, Y. Sun, YJ Suh, L. Shen, S. Yang, Smiling Women Pitching down : Audit des biais de genre représentationnels et présentationnels dans l'IA génératrice d'images. Journal of Computer-Mediated Communication : JCMC 29, zmad045 (2023) ; <https://doi.org/10.1093/jcmc/zmad045>.
- 497 Y. Wan, K.-W. Chang, Le PDG masculin et l'assistante féminine : évaluation et atténuation des préjugés sexistes dans Génération de texte en image de sujets doubles, arXiv [cs.CV] (2024) ; <http://arxiv.org/abs/2402.11089>.
- 498 A. Nielsen, A. Woemmel, « Inégalités invisibles : faire face à la discrimination fondée sur l'âge dans la recherche et les applications en apprentissage automatique » dans 2e atelier sur l'IA générative et le droit (2024) ; https://blog.genlaw.org/pdfs/genlaw_icml2024/50.pdf.
- 499 C. Harris, Atténuer les biais liés à l'âge dans les modèles d'IA de sélection des CV. Conférence internationale FLAIRS Actes 36 (2023) ; <https://doi.org/10.32473/flairs.36.133236>.
- 500 J. Stypinska, AI Ageism : une feuille de route essentielle pour l'étude de la discrimination et de l'exclusion liées à l'âge dans les environnements numériques Sociétés. IA & Société 38, 665–677 (2023) ; <https://doi.org/10.1007/s00146-022-01553-5>.
- 501 R. Naik, B. Nushi, « Les biais sociaux à travers le prisme de la génération de texte en image » dans les actes de la conférence 2023 Conférence AAAI/ACM sur l'IA, l'éthique et la société (AIES '23) (Association for Computing Machinery, New York, NY, États-Unis, 2023), pp. 786–808 ; <https://doi.org/10.1145/3600211.3604711>.
- 502* A. Tamkin, A. Askeel, L. Lovitt, E. Durmus, N. Joseph, S. Kravec, K. Nguyen, J. Kaplan, D. Ganguli, Evaluating and Atténuer la discrimination dans les décisions relatives aux modèles linguistiques, arXiv [cs.CL] (2023) ; <http://arxiv.org/abs/2312.03689>.

- 503 M. Kamruzzaman, Shovon, G. Kim, Enquête sur les biais subtils dans les LLM : âgisme, beauté, biais institutionnel et national dans les modèles génératifs (Association for Computational Linguistics, 2024) ; <https://doi.org/10.18653/v1/2024.findings-acl.530>.
- 504 CH Chu, S. Donato-Woodger, SS Khan, R. Nyrup, K. Leslie, A. Lyn, T. Shi, A. Bianchi, SA Rahimi, A. Grenier, Age-Biais connexes et intelligence artificielle : une revue de la littérature. *Humanities & Social Sciences Communications* 10, 1–17 (2023) ; <https://doi.org/10.1057/s41599-023-01999-y>.
- 505 T. Kamelski, D. Klinge, Intelligence artificielle générative et âgisme numérique : exploration de la construction de l'âge et Vieillesse par IA génératrice d'images (2024) ; <https://doi.org/10.31219/osf.io/p3sdj>.
- 506 KA Mack, R. Qadri, R. Denton, SK Kane, CL Bennett, « Ils ne se soucient que de nous montrer le fauteuil roulant » : représentation du handicap dans les modèles d'IA de conversion de texte en image » dans Actes de la conférence CHI sur les facteurs humains dans les systèmes informatiques (ACM, New York, NY, États-Unis, 2024) vol. 22, pp. 1–23 ; <https://doi.org/10.1145/3613904.3642166>.
- 507 PN Venkit, M. Srinath, S. Wilson, « Validisme automatisé : une exploration des biais explicites liés au handicap dans les sentiments Français et modèles d'analyse de toxicité » dans Actes du 3e atelier sur le traitement fiable du langage naturel (TrustNLP 2023), A. Ovalle, K.-W. Chang, N. Mehrabi, Y. Pruksachatkun, A. Galystan, J. Dhamala, A. Verma, T. Cao, A. Kumar, R. Gupta, éd. (Association for Computational Linguistics, Toronto, Canada, 2023), pp. 26–34 ; <https://doi.org/10.18653/v1/2023.trustnlp-1.3>.
- 508 K. Glazko, Y. Mohammed, B. Kosa, V. Potluri, J. Mankoff, « Identifier et améliorer les préjugés liés au handicap dans la sélection des CV basée sur le GPT » dans la conférence 2024 de l'ACM sur l'équité, la responsabilité et la transparence (ACM, New York, NY, États-Unis, 2024) ; <https://doi.org/10.1145/3630106.3658933>.
- 509 N. Shahin, L. Ismail, « ChatGPT, Let Us Chat Sign Language : Expériences, éléments architecturaux, défis et orientations de recherche » dans Symposium international 2023 sur les réseaux, les ordinateurs et les communications (ISNCC) (IEEE, 2023), pp. 1–7 ; <https://doi.org/10.1109/isncc58260.2023.10323974>.
- 510 S. Gueuwou, K. Takyi, M. Müller, MS Nyarko, R. Adade, R.-MOM Gyening, « AfriSign : traduction automatique pour les langues des signes africaines » dans 4e atelier sur le traitement du langage naturel africain (AfricaNLP 2023) (La Nouvelle-Orléans, LA, États-Unis, 2023) ; <https://openreview.net/forum?id=EHldk3J2xk>.
- 511 J. Hartmann, J. Schwenzow, M. Witte, L'idéologie politique de l'IA conversationnelle : preuves convergentes sur l'orientation pro-environnementale et libertaire de gauche de ChatGPT, arXiv [cs.CL] (2023) ; <http://arxiv.org/abs/2301.01768>.
- 512 F. Motoki, V. Pinho Neto, V. Rodrigues, Plus humain qu'humain : mesurer les biais politiques de ChatGPT. *Choix public* 198, 3-23 (2024) ; <https://doi.org/10.1007/s11127-023-01097-2>.
- 513 D. Rozado, Les biais politiques de ChatGPT. *Sciences sociales* (Bâle, Suisse) 12, 148 (2023) ; <https://doi.org/10.3390/socsci12030148>.
- 514 J. Rutinowski, S. Franke, J. Endendyk, I. Dormuth, M. Roidl, M. Pauly, La perception de soi et les préjugés politiques des ChatGPT. *Comportement humain et technologies émergentes* 2024, 1–9 (2024) ; <https://doi.org/10.1155/2024/7115633>.
- 515 M. Buyl, A. Rogiers, S. Noels, I. Dominguez-Catena, E. Heiter, R. Romero, I. Johary, A.-C. Mara, J. Lijffijt, T. De Bie, Les grands modèles de langage reflètent l'idéologie de leurs créateurs, arXiv [cs.CL] (2024) ; <http://arxiv.org/abs/2410.18417>.
- 516* T. Choudhary, Biais politique dans les modèles IA-langage : une analyse comparative de ChatGPT-4, Perplexity, Google Gemini et Claude, Techrxiv (2024) ; <https://doi.org/10.36227/techrxiv.172107441.12283354/v1>.
- 517 S. Feng, CY Park, Y. Liu, Y. Tsvetkov, Des données de pré-formation aux modèles linguistiques jusqu'aux tâches en aval : suivre les traces des préjugés politiques menant à des modèles NLP injustes (Association for Computational Linguistics, 2023) ; <https://doi.org/10.18653/v1/2023.acl-long.656>.
- 518 L. Rettenberger, M. Reischl, M. Schutera, Évaluation des préjugés politiques dans les grands modèles linguistiques, arXiv [cs.CL] (2024) ; <http://arxiv.org/abs/2405.13041>.
- 519 S. Fujimoto, K. Takemoto, Revisiter les biais politiques de ChatGPT. *Frontiers in Artificial Intelligence* 6, 1232003 (2023) ; <https://doi.org/10.3389/frai.2023.1232003>.
- 520 C. Walker, JC Timoneda, Identification des sources de biais idéologiques dans les modèles GPT grâce à la variation linguistique dans Sortie, arXiv [cs.CL] (2024) ; <http://arxiv.org/abs/2409.06043>.
- 521 T. Ceron, N. Falk, A. Barić, D. Nikolaev, S. Padó, Au-delà de la fragilité rapide : évaluation de la fiabilité et Cohérence des visions du monde politiques dans les LLM. *Transactions of the Association for Computational Linguistics* 12, 1378–1400 (2024) ; https://doi.org/10.1162/tacl_a_00710.
- 522 E. Perez, S. Ringer, K. Lukosiute, K. Nguyen, E. Chen, S. Heiner, C. Pettit, C. Olsson, S. Kundu, S. Kadavath, A. Jones, A. Chen, B. Mann, B. Israel, B. Seethor, C. McKinnon, C. Olah, ... J. Kaplan, « Découverte de comportements de modèles linguistiques à l'aide d'évaluations écrites de modèles » dans Conclusions de l'Association pour la linguistique computationnelle : ACL 2023, A. Rogers, J. Boyd-Graber, N. Okazaki, éd. (Association de linguistique computationnelle, Toronto, Canada, 2023), pp.

13387-13434 ; <https://doi.org/10.18653/v1/2023.findings-acl.847>.

- 523 J. Fisher, S. Feng, R. Aron, T. Richardson, Y. Choi, DW Fisher, J. Pan, Y. Tsvetkov, K. Reinecke, L'IA biaisée peut influencer la prise de décision politique, arXiv [cs.HC] (2024); <http://arxiv.org/abs/2410.06415>.
- 524 U. Messer, Comment les gens réagissent-ils aux préjugés politiques dans l'intelligence artificielle générative (IA) ? Les ordinateurs dans le comportement humain : humains artificiels, 100108 (2024) ; <https://doi.org/10.1016/j.chbah.2024.100108>.
- 525 Á. A. Cabrera, W. Epperson, F. Hohman, M. Kahng, J. Morgenstern, DH Chau, « FAIRVIS : Analyse visuelle pour « Découverte du biais intersectionnel dans l'apprentissage automatique » dans la conférence IEEE 2019 sur la science et la technologie de l'analyse visuelle (VAST) (2019), pp. 46–56 ; <https://doi.org/10.1109/VAST47406.2019.8986948>.
- 526 W. Guo, A. Caliskan, « Détection des biais intersectionnels émergents : les plongements de mots contextualisés contiennent un « Distribution of Human-like Biases » dans les actes de la conférence 2021 AAAI/ACM sur l'IA, l'éthique et la société (AIES '21) (Association for Computing Machinery, New York, NY, États-Unis, 2021), pp. 122–133 ; <https://doi.org/10.1145/3461702.3462536>.
- 527 IMS Lassen, M. Almasi, K. Enevoldsen, RD Kristensen-McLachlan, « Détection de l'intersectionnalité dans les modèles NER : une approche basée sur les données » dans Actes du 7e atelier conjoint SIGHUM sur la linguistique computationnelle pour le patrimoine culturel, les sciences sociales, les sciences humaines et la littérature, S. Degaetano-Ortlieb, A. Kazantseva, N. Reiter, S. Szpakowicz, éd. (Association pour la linguistique computationnelle, Dubrovnik, Croatie, 2023), pp. 116–127 ; <https://doi.org/10.18653/v1/2023.latechclfl-1.13>.
- 528 A. Ovalle, A. Subramonian, V. Gautam, G. Gee, K.-W. Chang, « Factorisation de la matrice de domination : un examen critique et une réimagination de l'intersectionnalité dans l'équité de l'IA » dans Actes de la conférence 2023 AAAI/ACM sur l'IA, l'éthique et la société (AIES '23) (Association for Computing Machinery, New York, NY, États-Unis, 2023), pp. 496–511 ; <https://doi.org/10.1145/3600211.3604705>.
- 529 K. Wilson, A. Caliskan, Genre, race et biais intersectionnel dans la sélection des CV via la recherche de modèles linguistiques, arXiv [cs.CY] (2024) ; <http://arxiv.org/abs/2407.20371>.
- 530 X. Fang, S. Che, M. Mao, H. Zhang, M. Zhao, X. Zhao, Biais du contenu généré par l'IA : un examen des actualités Produit par Large Language Models. Scientific Reports 14, 5224 (2024) ; <https://doi.org/10.1038/s41598-024-55686-2>.
- 531 H. An, C. Acquaye, C. Wang, Z. Li, R. Rudinger, « Les grands modèles linguistiques font-ils preuve de discrimination dans les décisions d'embauche en fonction de la race, de l'origine ethnique et du sexe ? » dans Actes de la 62e réunion annuelle de l'Association for Computational Linguistics (Volume 2 : Articles courts) (Association for Computational Linguistics, Stroudsburg, PA, États-Unis, 2024), pp. 386–397 ; <https://doi.org/10.18653/v1/2024.acl-short.37>.
- 532 R. Navigli, S. Conia, B. Ross, Biais dans les grands modèles linguistiques : origines, inventaire et discussion. J. Data and Information Quality 15, 1–21 (2023) ; <https://doi.org/10.1145/3597307>.
- 533 Y. Li, M. Du, R. Song, X. Wang, Y. Wang, Une enquête sur l'équité dans les grands modèles linguistiques, arXiv [cs.CL] (2023) ; <http://arxiv.org/abs/2308.10149>.
- 534* S. Mukherjee, A. Mitra, G. Jawahar, S. Agarwal, H. Palangi, A. Awadallah, Orca : Apprentissage progressif à partir de modèles complexes Explication des traces de GPT-4, arXiv [cs.CL] (2023) ; <http://arxiv.org/abs/2306.02707>.
- 535 E. Ferrara, Équité et biais en intelligence artificielle : un bref aperçu des sources, des impacts et des stratégies d'atténuation. Sci 6, 3 (2023) ; <https://doi.org/10.3390/sci6010003>.
- 536 SU Noble, Algorithmes d'oppression : comment les moteurs de recherche renforcent le racisme, NYU Press (2019) ; <https://nyupress.org/9781479837243/algorithms-of-oppression/>.
- 537 S. Lazar, A. Nelson, La sécurité de l'IA selon les conditions de qui ? Science 381, 138 (2023) ; <https://doi.org/10.1126/science.adi8982>.
- 538 RIJ Dobbe, TK Gilbert, Y. Mintz, « Choix difficiles en intelligence artificielle : faire face à l'incertitude normative « par le biais d'engagements sociotechniques (AIES '20) » dans Actes de la conférence AAAI/ACM sur l'IA, l'éthique et la société (Association for Computing Machinery, New York, NY, États-Unis, 2020), p. 242 ; <https://doi.org/10.1145/3375627.3375861>.
- 539 M. Shur-Ofry, La multiplicité comme principe de gouvernance de l'IA (2023) ; https://papers.ssrn.com/sol3/papers.cfm?abstract_id=4444354.
- 540 M. Sloane, E. Moss, O. Awomolo, L. Forlano, « La participation n'est pas une solution de conception pour l'apprentissage automatique » dans Actes de la 2e conférence de l'ACM sur l'équité et l'accès aux algorithmes, mécanismes et optimisation (EAAMO '22) (Association for Computing Machinery, New York, NY, États-Unis, 2022), pp. 1–6 ; <https://doi.org/10.1145/3551624.355285>.
- 541 H. Gonen, Y. Goldberg, « Du rouge à lèvres sur un cochon : les méthodes de débiasing dissimulent les préjugés sexistes systématiques dans Word « Incorporations mais ne les supprimez pas » dans les actes de l'atelier 2019 sur l'élargissement de la PNL, A. Axelrod, D. Yang, R. Cunha, S. Shaikh, Z. Waseem, éd. (Association de linguistique computationnelle, Florence, Italie, 2019), pp.

- 60-63 ; <https://aclanthology.org/W19-3621>.
- 542 J. Xiao, Z. Li, X. Xie, E. Getzen, C. Fang, Q. Long, WJ Su, Sur le biais algorithmique de l'alignement du langage volumineux Modèles avec RLHF : effondrement des préférences et régularisation correspondante, arXiv [stat.ML] (2024) ; <http://arxiv.org/abs/2405.16455>.
- 543 DY Kim, C. Wallraven, « Qualité des étiquettes dans AffectNet : résultats de la réannotation basée sur la foule » dans Lecture Notes in Computer Science (Springer International Publishing, Cham, 2022) Notes de cours en informatique, pp. 518–531 ; https://doi.org/10.1007/978-3-031-02444-3_39.
- 544 J. Ma, Y. Ushiku, M. Sagara, « L'effet de l'amélioration de la qualité des annotations sur les ensembles de données de détection d'objets : A Étude préliminaire » dans les ateliers de la conférence IEEE/CVF 2022 sur la vision par ordinateur et la reconnaissance de formes (CVPRW) (2022), pp. 4849–4858 ; <https://doi.org/10.1109/CVPRW56347.2022.00532>.
- 545 Z. Xu, K. Peng, L. Ding, D. Tao, X. Lu, Prenez soin de votre biais d'invite ! Enquête et atténuation du biais d'invite dans l'extraction de connaissances factuelles, arXiv [cs.CL] (2024) ; <http://arxiv.org/abs/2403.09963>.
- 546 H. Weerts, F. Pfisterer, M. Feurer, K. Eggensperger, E. Bergman, N. Awad, J. Vanschoren, M. Pechenizkiy, B. Bischl, F. Hutter, L'équité peut-elle être automatisée ? Lignes directrices et opportunités pour un AutoML soucieux de l'équité. The Journal of Artificial Intelligence Research 79, 639–677 (2024) ; <https://doi.org/10.1613/jair.1.14747>.
- 547 ID Raji, J. Buolamwini, « Audit exploitable : enquête sur l'impact de la dénonciation publique des performances biaisées Résultats des produits d'IA commerciaux » dans les actes de la conférence 2019 de l'AAAI/ACM sur l'IA, l'éthique et la société (ACM, New York, NY, États-Unis, 2019) ; <https://doi.org/10.1145/3306618.3314244>.
- 548 D. Zhang, P. Finckenberg-Broman, T. Hoang, S. Pan, Z. Xing, M. Staples, X. Xu, Le droit à l'oubli à l'ère des grands modèles linguistiques : implications, défis et solutions. IA et éthique (2024) ; <https://doi.org/10.1007/s43681-024-00573-9>.
- 549* A. Xiang, « Être vu » ou « mal vu » : tensions entre confidentialité et équité dans la vision par ordinateur. Harvard Journal of Law & Technology 36 (2022) ; https://papers.ssrn.com/sol3/papers.cfm?abstract_id=4068921.
- 550 J. Kleinberg, « Compromis inhérents à l'équité algorithmique » dans Résumés de l'ACM International 2018 Conférence sur la mesure et la modélisation des systèmes informatiques (SIGMETRICS '18) (Association for Computing Machinery, New York, NY, États-Unis, 2018), p. 40 ; <https://doi.org/10.1145/3219617.3219634>.
- 551 H. Nilforoshan, JD Gaebler, R. Shroff, S. Goel, « Conceptions causales de l'équité et leurs conséquences » dans Actes de la 39e Conférence internationale sur l'apprentissage automatique (ICML 2022) (PMLR, 2022) ; <https://proceedings.mlr.press/v162/nilforoshan22a.html>.
- 552 N. Konstantinov, CH Lampert, « Sur l'impossibilité d'un apprentissage soucieux de l'équité à partir de données corrompues » dans Atelier sur l'équité algorithmique à travers le prisme de la causalité et de la robustesse (AFCR 2021) (PMLR, virtuel, 2021) ; <https://proceedings.mlr.press/v171/konstantinov22a.html>.
- 553 A. Chouldechova, Prédiction équitable avec impact disparate : une étude des biais dans les instruments de prédiction de la récidive. Big Data 5, 153–163 (2017) ; <https://doi.org/10.1089/big.2016.0047>.
- 554 Q. Zhang, J. Liu, Z. Zhang, J. Wen, B. Mao, X. Yao, Atténuer l'injustice via un ensemble multi-objectifs évolutif Apprentissage. IEEE Transactions on Evolutionary Computation 27, 848–862 (2023) ; <https://doi.org/10.1109/TEVC.2022.3209544>.
- 555 M. Hardt, E. Price, E. Price, N. Srebro, « Égalité des chances dans l'apprentissage supervisé » dans 30e Conférence sur les systèmes de traitement de l'information neuronale (NIPS 2016) (Curran Associates, Inc., Barcelone, Espagne, 2016) vol. 29 ; https://proceedings.neurips.cc/paper_files/paper/2016/hash/9d2682367c3935defcb1f9e247a97c0d-Abstract.html.
- 556 M. Brcic, RV Yampolskiy, Résultats d'impossibilité en IA : une enquête. ACM Comput. Surv. 56, 1–24 (2023) ; <https://doi.org/10.1145/3603371>.
- 557 B. Green, Échapper à l'impossibilité de l'équité : de l'équité algorithmique formelle à l'équité algorithmique substantielle. Philosophie et Technologie 35, 90 (2022) ; <https://doi.org/10.1007/s13347-022-00584-6>.
- 558 A. Bell, L. Bynum, N. Drushchak, T. Zakharchenko, L. Rosenblatt, J. Stoyanovich, « La possibilité de l'équité : revisiter le théorème d'impossibilité dans la pratique » dans Actes de la conférence 2023 de l'ACM sur l'équité, la responsabilité et la transparence (FAccT '23) (Association for Computing Machinery, New York, NY, États-Unis, 2023), pp. 400–422 ; <https://doi.org/10.1145/3593013.3594007>.
- 559 KT Rodolfa, H. Lamba, R. Ghani, Observation empirique de compromis négligeables entre équité et précision dans l'apprentissage automatique pour les politiques publiques. Nature Machine Intelligence 3, 896–904 (2021) ; <https://doi.org/10.1038/s42256-021-00396-x>.
- 560 V. Hofmann, PR Kalluri, D. Jurafsky, S. King, Les préjugés dialectaux prédisent les décisions de l'IA sur le caractère, l'employabilité et la criminalité des personnes, arXiv [cs.CL] (2024) ; <http://arxiv.org/abs/2403.00742>.

- 561 RL Johnson, G. Pistilli, N. Menéndez-González, LDD Duran, E. Panai, J. Kalpokiene, DJ Bertulfo, Le fantôme dans la machine a un accent américain : conflit de valeurs dans GPT-3, arXiv [cs.CL] (2022) ; <http://arxiv.org/abs/2203.07785>.
- 562 E. Durmus, K. Nguyen, T. Liao, N. Schiefer, A. Askell, A. Bakhtin, C. Chen, Z. Hatfield-Dodds, D. Hernandez, N. Joseph, L. Lovitt, S. McCandlish, O. Sikder, A. Tamkin, J. Thamkul, J. Kaplan, J. Clark, D. Ganguli, « Vers la mesure de la représentation des opinions globales subjectives dans les modèles linguistiques » dans Première conférence sur la modélisation du langage (2024) ; <https://openreview.net/pdf?id=zl16jLb91v>.
- 563 Y. Wan, K.-W. Chang, Les hommes blancs dirigent, les femmes noires aident ? Analyse comparative des préjugés sociaux en matière d'agence linguistique LLM, arXiv [cs.CL] (2024) ; <http://arxiv.org/abs/2404.10508>.
- 564 B. Alkhamissi, M. Einokrashy, M. Alkhamissi, M. Diab, « Enquête sur l'alignement culturel des grands modèles linguistiques » dans Actes de la 62e réunion annuelle de l'Association for Computational Linguistics (Volume 1 : Long Papers) (Association for Computational Linguistics, Stroudsburg, PA, États-Unis, 2024), pp. 12404–12422 ; <https://doi.org/10.18653/v1/2024.acl-long.671>.
- 565 H. Yuan, Z. Che, S. Li, Y. Zhang, X. Hu, S. Luo, Le profil psychologique de haute dimension et le biais culturel de ChatGPT, arXiv [cs.CL] (2024) ; <http://arxiv.org/abs/2405.03387>.
- 566 R. Hada, S. Husain, V. Gumma, H. Diddee, A. Yadavalli, A. Seth, N. Kulkarni, U. Gadiraju, A. Vashistha, V. Seshadri, K. Bali, « Akal Badi Ya Bias : une étude exploratoire des préjugés sexistes dans la technologie de la langue hindi » dans la conférence 2024 de l'ACM sur l'équité, la responsabilité et la transparence (ACM, New York, NY, États-Unis, 2024) ; <https://doi.org/10.1145/3630106.3659017>.
- 567 MHJ Lee, JM Montgomery, CK Lai, « Les grands modèles linguistiques décrivent les groupes socialement subordonnés comme plus homogènes, ce qui est cohérent avec un biais observé chez les humains » dans la conférence 2024 de l'ACM sur l'équité, la responsabilité et la transparence (ACM, New York, NY, États-Unis, 2024) ; <https://doi.org/10.1145/3630106.3658975>.
- 568 C. Raj, A. Mukherjee, A. Caliskan, A. Anastasopoulos, Z. Zhu, Breaking Bias, Building Bridges: Évaluation et Atténuation des biais sociaux dans les LLM via l'hypothèse de contact. Actes de la conférence AAAI/ACM sur l'IA, l'éthique et la société 7, 1180–1189 (2024) ; <https://ojs.aaai.org/index.php/AIES/article/view/31715>.
- 569 D. Oba, M. Kaneko, D. Bollegala, « Suppression des préjugés sexistes en contexte pour les grands modèles linguistiques » dans Findings de l'Association pour la linguistique computationnelle : EAACL 2024 (2024), pp. 1722–1742 ; <https://aclanthology.org/2024.findings-eacl.121.pdf>.
- 570 Y. Reif, R. Schwartz, « Au-delà des performances : quantification et atténuation du biais d'étiquetage dans les LLM » dans Actes de la conférence 2024 du chapitre nord-américain de l'Association for Computational Linguistics : Human Language Technologies (Volume 1 : Long Papers) (Association for Computational Linguistics, Stroudsburg, PA, États-Unis, 2024), pp. 6784–6798 ; <https://doi.org/10.18653/v1/2024.naacl-long.378>.
- 571 M. Ribeiro, B. Malcorra, NB Mota, R. Wilkens, A. Villavicencio, LC Hubner, C. Rennó-Costa, A Methodology for Modèles linguistiques explicables à grande échelle avec gradients intégrés et analyse linguistique dans la classification de textes, arXiv [cs.CL] (2024) ; <http://arxiv.org/abs/2410.00250>.
- 572 L. Luo, Y.-F. Li, R. Haf, S. Pan, « Raisonement sur les graphiques : raisonnement fidèle et interprétable à partir d'un grand modèle de langage » dans 12e Conférence internationale sur les représentations d'apprentissage (2023) ; <https://openreview.net/pdf?id=ZGNWW7xZ6Q>.
- 573 S. Pan, L. Luo, Y. Wang, C. Chen, J. Wang, X. Wu, Unifier les grands modèles de langage et les graphes de connaissances : une feuille de route. IEEE Transactions on Knowledge and Data Engineering 36, 3580–3599 (2024) ; <https://doi.org/10.1109/tkde.2024.3352100>.
- 574 SA Friedler, C. Scheidegger, S. Venkatasubramanian, La (im)possibilité d'équité : différents systèmes de valeurs nécessitent différents mécanismes pour une prise de décision équitable. Communications de l'ACM 64, 136–143 (2021) ; <https://doi.org/10.1145/3433949>.
- 575 J. Banja, JW Gichoya, N. Martinez-Martin, LA Waller, GD Clifford, L'équité comme réflexion après coup : une perspective américaine sur l'équité dans les collaborations entre développeurs de modèles et utilisateurs cliniciens. PLOS Digital Health 2, e0000386 (2023) ; <https://doi.org/10.1371/journal.pdig.0000386>.
- 576 NA Saxena, K. Huang, E. DeFilippis, G. Radanovic, DC Parkes, Y. Liu, « Comment se comportent les définitions d'équité ? « Examen des attitudes du public à l'égard des définitions algorithmiques de l'équité » dans les actes de la conférence 2019 de l'AAAI/ACM sur l'IA, l'éthique et la société (AIES '19) (Association for Computing Machinery, New York, NY, États-Unis, 2019), pp. 99–106 ; <https://doi.org/10.1145/3306618.3314248>.
- 577 W. Fleisher, « Qu'est-ce qui est juste dans l'équité individuelle ? » dans Actes de la conférence 2021 de l'AAAI/ACM sur l'IA, l'éthique et la société (AIES '21) (Association for Computing Machinery, New York, NY, États-Unis, 2021), pp. 480–490 ; <https://doi.org/10.1145/3461702.3462621>.
- 578 AM Turing, Intelligent Machinery, A Heretical Theory*. Philosophia Mathematica. Série III 4, 256–260 (1996) ;

- <https://doi.org/10.1093/philmat/4.3.256>.
- 579 IJ Good, « Spéculations concernant la première machine ultra-intelligente » dans *Advances in Computers*, FL Alt, M. Rubinoff, éd. (Elsevier, 1966) vol. 6, pp. [https://doi.org/10.1016/S0065-2458\(08\)60418-0](https://doi.org/10.1016/S0065-2458(08)60418-0).
- 580 N. Wiener, Quelques conséquences morales et techniques de l'automatisation. *Science* 131, 1355–1358 (1960) ; <https://doi.org/10.1126/science.131.3410.1355>.
- 581 SM Omohundro, « Les moteurs de base de l'IA » dans les actes de la conférence de 2008 sur l'intelligence artificielle générale 2008 : Actes de la première conférence de l'AGI (IOS Press, NLD, 2008), pp. 483–492 ; <https://dl.acm.org/doi/10.5555/1566174.1566226>.
- 582 N. Bostrom, *MM Cirkovic, Risques catastrophiques mondiaux* (Oxford University Press, Londres, Angleterre, 2011) ; <https://academic.oup.com/book/40615>.
- 583 S. Russell, P. Norvig, *Intelligence artificielle* (Pearson, Upper Saddle River, NJ, éd. 3, 2009) ; https://people.engr.tamu.edu/guni/csce421/files/AI_Russell_Norvig.pdf.
- 584 N. Bostrom, *Superintelligence : chemins, dangers, stratégies* (Oxford University Press, Londres, Angleterre, 2014) ; <https://global.oup.com/academic/product/superintelligence-9780198739838?cc=mx&lang=en&>.
- 585 SJ Russell, *Human Compatible : L'intelligence artificielle et le problème du contrôle* (Penguin Books, Harlow, Angleterre, 2020) ; <https://www.penguin.co.uk/books/307948/human-compatible-by-russell-stuart/9780141987507> .
- 586 Centre pour la sécurité de l'IA, Déclaration sur les risques liés à l'IA : Les experts en IA et les personnalités publiques expriment leur inquiétude face aux risques liés à l'IA (2024) ; <https://www.safe.ai/work/statement-on-ai-risk>.
- 587 Y. Bengio, Déclaration écrite du professeur Yoshua Bengio devant le Forum du Sénat américain sur l'intelligence artificielle concernant le risque, l'alignement et la protection contre les scénarios apocalyptiques. (2023) ; <https://www.schumer.senate.gov/imo/media/doc/Yoshua%20Benigo%20-%20Statement.pdf>.
- 588 K. Grace, H. Stewart, JF Sandkühler, S. Thomas, B. Weinstein-Raun, J. Brauner, Des milliers d'auteurs d'IA sur le L'avenir de l'IA, arXiv [cs.CY] (2024) ; <http://arxiv.org/abs/2401.02843>.
- 589 A. Critch, S. Russell, TASRA : Une taxonomie et une analyse des risques à l'échelle sociétale liés à l'IA, arXiv [cs.AI] (2023) ; <http://arxiv.org/abs/2306.06924>.
- 590 K. Goddard, A. Roudsari, JC Wyatt, Biais d'automatisation : une revue systématique de la fréquence, des médiateurs d'effets et des atténuateurs. *Journal de l'American Medical Informatics Association : JAMIA* 19, 121–127 (2012) ; <https://doi.org/10.1136/amiajnl-2011-000089>.
- 591 M. Chugunova, D. Sele, Nous et cela : une revue interdisciplinaire des preuves expérimentales sur la façon dont les humains interagissent avec les machines. *Journal of Behavioral and Experimental Economics* 99, 101897 (2022) ; <https://doi.org/10.1016/j.socec.2022.101897>.
- 592 A. Kasirzadeh, Deux types de risques existentiels liés à l'IA : décisif et cumulatif, arXiv [cs.CY] (2024) ; <http://arxiv.org/abs/2401.07836>.
- 593 M. Kinniment, LJK Sato, H. Du, B. Goodrich, M. Hasin, L. Chan, LH Miles, TR Lin, H. Wijk, J. Burget, A. Ho, E. Barnes, P. Christiano, Évaluation des agents de modèles de langage sur des tâches autonomes réalistes, arXiv [cs.CL] (2023) ; https://evals.alignment.org/Evaluating_LMAs_Realistic_Tasks.pdf.
- 594* OpenAI, « Cadre de préparation (bêta) » (OpenAI, 2023) ; <https://cdn.openai.com/openai-preparedness-framework-beta.pdf> .
- 595* Anthropic, Politique de mise à l'échelle responsable d'Anthropic, version 1.0. (2023) ; <https://www-cdn.anthropic.com/1adf000c8f675958c2ee23805d91aaade1cd4613/responsible-scaling-policy.pdf>.
- 596* Google DeepMind, Cadre de sécurité Frontier Version 1.0. (2024) ; <https://storage.googleapis.com/deepmind-media/DeepMind.com/Blog/introduction-du-cadre-de-sécurité-frontière/rapport-technique-fsf.pdf>.
- 597 T. Hagendorff, Les capacités de tromperie émergent dans les grands modèles linguistiques. *Actes de l'Académie nationale des sciences des États-Unis d'Amérique* 121, e2317967121 (2024) ; <https://doi.org/10.1073/pnas.2317967121>.
- 598* E. Hubinger, C. Denison, J. Mu, M. Lambert, M. Tong, M. MacDiarmid, T. Lanham, DM Ziegler, T. Maxwell, N. Cheng, A. Jermyn, A. Askeel, A. Radhakrishnan, C. Anil, D. Duvenaud, D. Ganguli, F. Barez, ... E. Perez, Agents dormants : former des LLM trompeurs qui. Persister grâce à la formation à la sécurité, arXiv [cs.CR] (2024) ; <http://arxiv.org/abs/2401.05566>.
- 599* C. Denison, M. MacDiarmid, F. Barez, D. Duvenaud, S. Kravec, S. Marks, N. Schiefer, R. Soklaski, A. Tamkin, J. Kaplan, B. Shlegeris, SR Bowman, E. Perez, E. Hubinger, De la flagornerie au subterfuge : étude de la falsification des récompenses dans les grands modèles linguistiques, arXiv [cs.AI] (2024) ; <http://arxiv.org/abs/2406.10162>.
- 600 S. Kapoor, B. Stroebel, ZS Siegel, N. Nadgir, A. Narayanan, AI Agents That Matter, arXiv [cs.LG] (2024) ;

- <http://arxiv.org/abs/2407.01502>.
- 601 R. Shiffrin, M. Mitchell, Sonder la psychologie des modèles d'IA. Actes de l'Académie nationale des sciences de États-Unis d'Amérique 120, e2300963120 (2023) ; <https://doi.org/10.1073/pnas.2300963120>.
- 602 D. Hendrycks, M. Mazeika, T. Woodside, Un aperçu des risques catastrophiques liés à l'IA, arXiv [cs.CY] (2023) ; <http://arxiv.org/abs/2306.12001>.
- 603 J. Lehman, J. Clune, D. Misevic, C. Adami, L. Altenberg, J. Beaulieu, PJ Bentley, S. Bernard, G. Beslon, DM Bryson, N. Cheney, P. Chrabaszcz, A. Cully, S. Doncieux, FC Dyer, KO Ellefsen, R. Feldt, ... J. Yosinski, La créativité surprenante de l'évolution numérique : un recueil d'anecdotes des communautés de recherche sur le calcul évolutionnaire et la vie artificielle. *Artificial Life* 26, 274–306 (2020) ; https://doi.org/10.1162/artl_a_00319.
- 604 J. Skalse, NHR Howe, D. Krashennikov, D. Krueger, Définition et caractérisation du piratage de récompense, arXiv [cs.LG] (2022) ; <http://arxiv.org/abs/2209.13085>.
- 605 R. Ngo, L. Chan, S. Mindermann, « Le problème d'alignement du point de vue de l'apprentissage profond » dans The 12th Conférence internationale sur les représentations de l'apprentissage (ICLR 2024) (Vienne, Autriche, 2023) ; <https://openreview.net/forum?id=fh8EYKFKns>.
- 606 J. Ji, T. Qiu, B. Chen, B. Zhang, H. Lou, K. Wang, Y. Duan, Z. He, J. Zhou, Z. Zhang, F. Zeng, KY Ng, J. Dai, X. Pan, A. O'Gara, Y. Lei, H. Xu, ... W. Gao, Alignement de l'IA : une étude approfondie, arXiv [cs.AI] (2023) ; <http://arxiv.org/abs/2310.19852>.
- 607 A. Pan, K. Bhatia, J. Steinhardt, « Les effets de la mauvaise spécification des récompenses : cartographie et atténuation des modèles mal alignés » dans la 10e Conférence internationale sur les représentations de l'apprentissage (ICLR 2022) (virtuelle, 2021) ; <https://openreview.net/forum?id=JYtwGwL7ye>.
- 608 J. Wen, R. Zhong, A. Khan, E. Perez, J. Steinhardt, M. Huang, SR Bowman, H. He, S. Feng, Modèles de langage Apprendre à Tromper les humains via RLHF, arXiv [cs.CL] (2024) ; <http://arxiv.org/abs/2409.12822>.
- 609* SR Bowman, J. Hyun, E. Perez, E. Chen, C. Pettit, S. Heiner, K. Lukošiuūtė, A. Askell, A. Jones, A. Chen, A. Goldie, A. Mirhoseini, C. McKinnon, C. Olah, D. Amodei, D. Amodei, D. Drain, ... J. Kaplan, Mesurer les progrès en matière de surveillance évolutive pour les grands modèles linguistiques, arXiv [cs.HC] (2022) ; <http://arxiv.org/abs/2211.03540>.
- 610* P. Christiano, B. Shlegeris, Dario, Amodei, Superviser les apprenants forts en amplifiant les experts faibles, arXiv [cs.LG] (2018) ; <http://arxiv.org/abs/1810.08575>.
- 611* G. Irving, P. Christiano, D. Amodei, « AI Safety via Debate » (OpenAI, 2018) ; <http://arxiv.org/abs/1805.00899>.
- 612* J. Leike, D. Krueger, T. Everitt, M. Martic, V. Maini, S. Legg, Alignement évolutif des agents via la modélisation des récompenses : A Direction de la recherche, arXiv [cs.LG] (2018) ; <http://arxiv.org/abs/1811.07871>.
- 613* J. Wu, L. Ouyang, DM Ziegler, N. Stiennon, R. Lowe, J. Leike, P. Christiano, Résumer récursivement des livres avec Rétroaction humaine, arXiv [cs.CL] (2021) ; <http://arxiv.org/abs/2109.10862>.
- 614* W. Saunders, C. Yeh, J. Wu, S. Bills, L. Ouyang, J. Ward, J. Leike, Modèles autocritiques pour aider les évaluateurs humains, arXiv [cs.CL] (2022) ; <http://arxiv.org/abs/2206.05802>.
- 615* A. Khan, J. Hughes, D. Valentine, L. Ruis, K. Sachan, A. Radhakrishnan, E. Grefenstette, SR Bowman, T. Rocktäschel, E. Perez, Débattre avec des LLM plus convaincants conduit à des réponses plus véridiques, arXiv [cs.AI] (2024) ; <http://arxiv.org/abs/2402.06782>.
- 616 LLD Langosco, J. Koch, LD Sharkey, J. Pfau, D. Krueger, « Généralisation erronée des objectifs dans le renforcement profond « Learning » dans les actes de la 39e Conférence internationale sur l'apprentissage automatique (PMLR, 2022) vol. 162, pp. 12004–12019 ; <https://proceedings.mlr.press/v162/langosco22a.html>.
- 617* R. Shah, V. Varma, R. Kumar, M. Phuong, V. Krakovna, J. Uesato, Z. Kenton, Généralisation erronée des objectifs : pourquoi des spécifications correctes ne suffisent pas pour des objectifs corrects, arXiv [cs.LG] (2022) ; <http://arxiv.org/abs/2210.01790>.
- 618 HNE Barj, T. Sautory, Apprentissage par renforcement à partir du feedback LLM pour contrer la généralisation erronée des objectifs, arXiv [cs.LG] (2024) ; <http://arxiv.org/abs/2401.07181>.
- 619 D. Hendrycks, X. Liu, E. Wallace, A. Dziedzic, R. Krishnan, D. Song, « Les transformateurs préentraînés améliorent la robustesse hors distribution » dans Actes de la 58e réunion annuelle de l'Association for Computational Linguistics (ACL 2020), D. Jurafsky, J. Chai, N. Schluter, J. Tetreault, éd. (Association for Computational Linguistics, en ligne, 2020), pp. 2744–2751 ; <https://doi.org/10.18653/v1/2020.acl-main.244>.
- 620 L. Berglund, AC Stickland, M. Balesni, M. Kaufmann, M. Tong, T. Korbak, D. Kokotajlo, O. Evans, extrait de Contexte : Sur la mesure de la connaissance de la situation dans les LLM, arXiv [cs.CL] (2023) ; <http://arxiv.org/abs/2309.00667>.
- 621 R. Laine, B. Chughtai, J. Betley, K. Hariharan, M. Balesni, J. Scheurer, M. Hobbhahn, A. Meinke, O. Evans, « Moi, moi-même et l'IA : l'ensemble de données de connaissance de la situation (SAD) pour les LLM » dans 38e conférence sur les ensembles de données et les repères des systèmes de traitement de l'information neuronale (2024) ; <https://openreview.net/forum?id=UnWhcplyUC>.

- 622 C. Schwab, L. Huber, Obéir ou ne pas obéir ? Les chiens (*Canis Familiaris*) se comportent différemment en réponse aux états d'attention de leurs propriétaires. *Journal of Comparative Psychology* (Washington, DC : 1983) 120, 169–175 (2006) ; <https://doi.org/10.1037/0735-7036.120.3.169>.
- 623* V. Krakovna, J. Kramar, La recherche de pouvoir peut être probable et prédictive pour les agents entraînés, arXiv [cs.AI] (2023) ; <http://arxiv.org/abs/2304.06528>.
- 624 A. Turner, L. Smith, R. Shah, A. Critch, P. Tadepalli, « Les politiques optimales tendent à rechercher le pouvoir » dans la 35e Conférence sur Systèmes de traitement de l'information neuronale (NeurIPS 2021) (Curran Associates, Inc., Virtuel, 2021) vol. 34 ; <https://proceedings.neurips.cc/paper/2021/hash/c26820b8a4c1b3c2aa868d6d57e14a79-Abstract.html>.
- 625 A. Turner, P. Tadepalli, « Les décideurs paramétriquement reciblables ont tendance à rechercher le pouvoir » dans *Advances in Neural Information Processing Systems (NeurIPS 2022) Main Conference Track* (La Nouvelle-Orléans, LA, États-Unis, 2022) vol. abs/2206.13477 ; <https://doi.org/10.48550/arXiv.2206.13477>.
- 626 MK Cohen, M. Hutter, MA Osborne, Des agents artificiels avancés interviennent dans la fourniture de récompenses. *AI Magazine* 43, 282–293 (2022) ; <https://doi.org/10.1002/aaai.12064>.
- 627 S. Zhuang, D. Hadfield-Menell, « Conséquences d'une IA mal alignée » dans *Advances in Neural Information Processing Systems (NeurIPS 2020)* (Curran Associates, Inc., 2020) vol. 33, pp. 15763–15773 ; <https://proceedings.neurips.cc/paper/2020/hash/b607ba543ad05417b8507ee86c54fcb7-Abstract.html>.
- 628 E. Hubinger, C. van Merwijk, V. Mikulik, J. Skalse, S. Garrabrant, Risques liés à l'optimisation apprise dans les systèmes d'apprentissage automatique avancés, arXiv [cs.AI] (2019) ; <http://arxiv.org/abs/1906.01820>.
- 629 J. Carlsmith, Scheming AIs : Les IA simuleront-elles l'alignement pendant l'entraînement afin d'obtenir du pouvoir ?, arXiv [cs.CY] (2023) ; <http://arxiv.org/abs/2311.08379>.
- 630* R. Grosse, J. Bae, C. Anil, N. Elhage, A. Tamkin, A. Tajdini, B. Steiner, D. Li, E. Durmus, E. Perez, E. Hubinger, K. Lukošiūtė, K. Nguyen, N. Joseph, S. McCandlish, J. Kaplan, SR Bowman, Étude de la généralisation de grands modèles de langage avec des fonctions d'influence, arXiv [cs.LG] (2023) ; <http://arxiv.org/abs/2308.03296>.
- 631 S. Im, Y. Li, Sur la généralisation de l'apprentissage des préférences avec DPO, arXiv [cs.LG] (2024) ; <http://arxiv.org/abs/2408.03459>.
- 632 A. Pan, JS Chan, A. Zou, N. Li, S. Basart, T. Woodside, H. Zhang, S. Emmons, D. Hendrycks, « Les récompenses justifient-elles les moyens ? Mesure des compromis entre récompenses et comportement éthique dans le référentiel MACHIAVELLI » dans *Actes de la 40e Conférence internationale sur l'apprentissage automatique (ICML'23) (JMLR, Honolulu, Hawaï, États-Unis, 2023)* vol. 202, pp. 26837–26867.
- 633 L. Dung, L'argument en faveur d'une déresponsabilisation humaine à court terme par l'IA. *AI & Society*, 1–14 (2024) ; <https://doi.org/10.1007/s00146-024-01930-2>.
- 634 PJ Denning, La science de l'informatique : le ver Internet. *American Scientist* 77, 126–128 (1989) ; <http://www.jstor.org/stable/27855650>.
- 635 D. Hendrycks, La sélection naturelle favorise les IA par rapport aux humains, arXiv [cs.CY] (2023) ; <http://arxiv.org/abs/2303.16200>.
- 636 UK AI Safety Institute, Faire progresser le domaine de la sécurité systémique de l'IA : subventions ouvertes (2024) ; <https://www.aisi.gov.uk/work/advancing-the-field-of-systemic-ai-safety-grants-open>.
- 637* T. Eloundou, S. Manning, P. Mishkin, D. Rock, Les GPT sont des GPT : le potentiel d'impact des LLM sur le marché du travail. *Science* 384, 1306-1308 (2024) ; <https://doi.org/10.1126/science.adj0998>.
- 638 B. Lou, H. Sun, T. Sun, GPTs et marchés du travail dans l'économie en développement : preuves de la Chine, SSRN [préimpression] (2023) ; <https://doi.org/10.2139/ssrn.4426461>.
- 639 P. Gmyrek, J. Berg, D. Bescond, IA générative et emplois : une analyse globale des effets potentiels sur la quantité d'emplois et qualité (Organisation internationale du travail, Genève, 2023) ; <https://doi.org/10.54394/fhem8239>.
- 640 M. Cazzaniga, F. Jaumotte, L. Li, G. Melina, AJ Panton, C. Pizzinelli, EJ Rockall, MM Tavares, « Gen-AI : l'intelligence artificielle et l'avenir du travail » (SDN/2024/001, Fonds monétaire international, 2024) ; <https://www.imf.org/en/Publications/Staff-Discussion-Notes/Issues/2024/01/14/Gen-AI-Artificial-Intelligence-and-the-Future-of-Work-542379>.
- 641 D. Acemoglu, P. Restrepo, Automatisation et nouvelles tâches : comment la technologie remplace et rétablit le travail. *The Journal of Economic Perspectives : A Journal of the American Economic Association* 33, 3–30 (2019) ; <https://doi.org/10.1257/jep.33.2.3>.
- 642 D. Acemoglu, D. Autor, « Compétences, tâches et technologies : implications pour l'emploi et les revenus* » dans *Manuel de l'économie du travail*, D. Card, O. Ashenfelter, éd. (Elsevier, 2011) vol. 4, pp. 1043–1171 ; [https://doi.org/10.1016/S0169-7218\(11\)02410-5](https://doi.org/10.1016/S0169-7218(11)02410-5).
- 643 P. Restrepo, « Automatisation : théorie, preuves et perspectives » (w31910, National Bureau of Economic Research, 2023) ; <https://doi.org/10.3386/w31910>.

- 644 D. Autor, C. Chin, A. Salomons, B. Seegmiller, « Nouvelles frontières : origines et contenu des nouvelles œuvres, 1940-2018 » (30389, Bureau national de recherche économique, 2022) ; <https://doi.org/10.3386/w30389>.
- 645 X. Hui, O. Reshef, L. Zhou, « Les effets à court terme de l'intelligence artificielle générative sur l'emploi : preuves « d'un marché du travail en ligne » (10601, document de travail CESifo, 2023) ; <https://www.econstor.eu/handle/10419/279352>.
- 646 A. Korinek, D. Suh, « Scénarios pour la transition vers l'AGI » (32255, National Bureau of Economic Research, 2024) ; <https://doi.org/10.3386/w32255>.
- 647 A. Korinek, Planification de scénarios pour un avenir d'IA(G). Magazine Finance et Développement (2023) ; <https://www.imf.org/en/Publications/fandd/issues/2023/12/Scenario-Planning-for-an-AGI-future-Anton-korinek>.
- 648 D. Acemoglu, « La macroéconomie simple de l'IA » (w32487, National Bureau of Economic Research, 2024) ; <https://doi.org/10.3386/w32487>.
- 649 B. Romera-Paredes, M. Barekatin, A. Novikov, M. Balog, MP Kumar, E. Dupont, FJR Ruiz, JS Ellenberg, P. Wang, O. Fawzi, P. Kohli, A. Fawzi, Découvertes mathématiques issues de la recherche de programmes avec de grands modèles de langage. Nature 625, 468–475 (2024) ; <https://doi.org/10.1038/s41586-023-06924-6>.
- 650 Y. Li, D. Choi, J. Chung, N. Kushman, J. Schrittwieser, R. Leblond, T. Eccles, J. Keeling, F. Gimeno, A. Dal Lago, T. Hubert, P. Choy, C. de Masson d'Autume, I. Babuschkin, X. Chen, P.-S. Huang, J. Welbl, ... O. Vinyals, Génération de code de niveau compétition avec AlphaCode. Science (New York, NY) 378, 1092-1097 (2022) ; <https://doi.org/10.1126/science.abq1158>.
- 651 S. Noy, W. Zhang, Preuves expérimentales des effets de l'intelligence artificielle générative sur la productivité. Science (New York, New York) 381, 187-192 (2023) ; <https://doi.org/10.1126/science.adh2586>.
- 652 D. Susskind, Un monde sans travail : technologie, automatisation et comment nous devrions réagir (Metropolitan Books, 2020) ; <https://www.danielsusskind.com/a-world-without-work>.
- 653 A. Korinek, M. Juelfs, « Se préparer à l'avenir (inexistant ?) du travail » (w30172, National Bureau of Economic Research, 2022) ; <https://doi.org/10.3386/w30172>.
- 654 A. Korinek, « Défis de politique économique à l'ère de l'IA » (w32980, National Bureau of Economic Research, 2024) ; <https://doi.org/10.3386/w32980>.
- 655* A. McAfee, « Généralement plus rapide : l'impact économique de l'IA générative » (Google, 2024) ; https://policycommons.net/artifacts/12281693/generally_faster_-_the_economic_impact_of_generative_ai/.
- 656 A. Agrawal, J. Gans, A. Goldfarb, « Adoption de l'IA et changement à l'échelle du système » (w28811, National Bureau of Economic Research, 2021) ; <https://doi.org/10.3386/w28811>.
- 657 J. Feigenbaum, DP Gross, Obstacles organisationnels et économiques à l'automatisation : un récit édifiant d'AT&T au XXe siècle. Management Science (2024) ; <https://doi.org/10.1287/mnsc.2022.01760>.
- 658 M. Svanberg, W. Li, M. Fleming, B. Goehring, N. Thompson, Au-delà de l'exposition à l'IA : quelles tâches sont rentables automatiser avec la vision par ordinateur ?, SSRN [préimpression] (2024) ; <https://doi.org/10.2139/ssrn.4700751>.
- 659 V. Magesh, F. Surani, M. Dahl, M. Suzgun, CD Manning, DE Ho, Hallucination-Free? Évaluation de la fiabilité des principaux outils de recherche juridique basés sur l'IA, arXiv [cs.CL] (2024) ; <http://arxiv.org/abs/2405.20362>.
- 660 E. Erdil, T. Besiroglu, Croissance explosive grâce à l'automatisation de l'IA : examen des arguments, arXiv [econ.GN] (2023) ; <https://epoch.ai/blog/explosive-growth-from-ai-a-review-of-the-arguments>.
- 661 A. Bick, A. Blandin, D. Deming, « L'adoption rapide de l'IA générative » (w32966, National Bureau of Economic Research, 2024) ; <https://doi.org/10.3386/w32966>.
- 662 E. Brynjolfsson, D. Li, L. Raymond, « L'IA générative au travail » (w31161, National Bureau of Economic Research, 2023) ; <https://doi.org/10.3386/w31161>.
- 663 D. Acemoglu, P. Restrepo, La course entre l'homme et la machine : implications de la technologie pour la croissance, Factor Actions et emploi. American Economic Review 108, 1488–1542 (2018) ; <https://doi.org/10.1257/aer.20160696>.
- 664 AK Agrawal, JS Gans, A. Goldfarb, « La transformation de Turing : intelligence artificielle, augmentation de l'intelligence et primes de compétences » (31767, National Bureau of Economic Research, 2023) ; <https://doi.org/10.3386/w31767>.
- 665 E. Felten, M. Raj, R. Seamans, Comment les modélisateurs de langage comme ChatGPT affecteront-ils les professions et les industries ?, arXiv [econ.GN] (2023) ; <http://arxiv.org/abs/2303.01157>.
- 666 EW Felten, M. Raj, R. Seamans, Hétérogénéité professionnelle dans l'exposition à l'IA générative, SSRN [préimpression] (2023) ; <https://doi.org/10.2139/ssrn.4414065>.
- 667 F. Dell'Acqua, E. McFowland III, ER Mollick, H. Lifshitz-Assaf, K. Kellogg, S. Rajendran, L. Kraymer, F. Candelon, KR

- Lakhani, « Naviguer dans la frontière technologique déshabillée : preuves expérimentales sur le terrain des effets de l'IA sur la productivité et la qualité des travailleurs du savoir » (24-013, Harvard Business School, 2023) ; https://www.hbs.edu/ris/Publication%20Files/24-013_d9b45b68-9e74-42d6-a1c6-c72fb70c7282.pdf.
- 668 JH Choi, A. Monahan, DB Schwarcz, L'exercice du droit à l'ère de l'intelligence artificielle, SSRN [préimpression] (2023) ; <https://doi.org/10.2139/ssrn.4626276>.
- 669 K. Bonney, C. Breaux, C. Buffington, E. Dinlersoz, L. Foster, N. Goldschlag, J. Haltiwanger, Z. Kroff, K. Savage, « Suivi de l'utilisation de l'IA par les entreprises en temps réel : un aperçu de l'enquête sur les tendances et les perspectives des entreprises » (w32319, National Bureau of Economic Research, 2024) ; <https://doi.org/10.3386/w32319>.
- 670 A. Korinek, L'essor des agents artificiellement intelligents (2019) ; https://drive.google.com/file/d/16y5UmeTOv5YB9E5ms_ce7WiYnFMA17J/view.
- 671 A. Chan, R. Salganik, A. Markelius, C. Pang, N. Rajkumar, D. Krasheninnikov, L. Langosco, Z. He, Y. Duan, M. Carroll, M. Lin, A. Mayhew, K. Collins, M. Molamohammadi, J. Burden, W. Zhao, S. Rismani, ... T. Maharaj, « Les préjugés causés par des systèmes algorithmiques de plus en plus agencés » dans Actes de la conférence 2023 de l'ACM sur l'équité, la responsabilité et la transparence (FAccT '23) (Association for Computing Machinery, New York, NY, États-Unis, 2023), pp. 651–666 ; <https://doi.org/10.1145/3593013.3594033>.
- 672 METR, Détails sur l'évaluation préliminaire du GPT-4o par METR, Ressources d'évaluation de l'autonomie du METR (2024) ; <https://metr.github.io/autonomy-evals-guide/gpt-4o-report/>.
- 673 Y. Shavit, S. Agarwal, M. Brundage, SAC O'Keefe, R. Campbell, T. Lee, P. Mishkin, T. Eloundou, A. Hickey, K. Slama, L. Ahmad, P. McMillan, A. Beutel, A. Passos, DG Robinson, Pratiques pour la gouvernance des systèmes d'IA agencés. Document de recherche, OpenAI (2023) ; <https://cdn.openai.com/papers/practices-for-governing-agentic-ai-systems.pdf>.
- 674 D. Hyslop, W. Townsend, « Les impacts à long terme des suppressions d'emplois sur les résultats du marché du travail » (Motu Recherche économique et politique publique, 2017) ; <https://www.motu.nz/our-research/population-and-labour/individual-and-group-outcomes/the-longer-term-impacts-of-job-displacement-on-labour-market-outcomes/>.
- 675 SC Dixon, DC Maré, « Les coûts de la perte d'emploi involontaire : impacts sur l'emploi et les revenus des travailleurs » (Recherche économique et politique publique Motu, 2013) ; <https://doi.org/10.2139/ssrn.2247198>.
- 676 D. Hamermesh, « Que savons-nous du déplacement des travailleurs aux États-Unis ? » (National Bureau of Economic Research, 1987) ; <https://doi.org/10.3386/w2402>.
- 677 LS Jacobson, RJ LaLonde, DG Sullivan, Pertes de revenus des travailleurs licenciés. *The American Economic Review* 83, 685–709 (1993) ; <http://www.jstor.org/stable/2117574>.
- 678 T. Von Wachter, J. Song, J. Manchester, Pertes de revenus à long terme dues aux licenciements massifs pendant la récession de 1982 : une analyse utilisant les données administratives américaines de 1974 à 2004 (2009) ; http://www.econ.ucla.edu/tvwachter/papers/mass_layoffs_1982.pdf.
- 679 J. Barnette, A. Michaud, Cicatrices salariales et théorie du capital humain (2017) ; <https://ammichau.github.io/papers/JBAMWageScar.pdf>.
- 680 D. Sullivan, T. von Wachter, Déplacements d'emplois et mortalité : une analyse à partir de données administratives*. *The Quarterly Journal of Economics* 124, 1265–1306 (2009) ; <https://doi.org/10.1162/qjec.2009.124.3.1265>.
- 681 SA Burgard, JE Brand, JS House, Vers une meilleure estimation de l'effet de la perte d'emploi sur la santé. *Journal of Health and Social Behavior* 48, 369–384 (2007) ; <http://www.jstor.org/stable/27638722>.
- 682 M. Browning, E. Heinesen, Effet de la perte d'emploi due à la fermeture d'une usine sur la mortalité et l'hospitalisation. *Journal of Health Economics* 31, 599–616 (2012) ; <https://doi.org/10.1016/j.jhealeco.2012.03.001>.
- 683 K. Telle, M. Votruba, Perte d'emploi des parents et résultats scolaires des enfants. *Revue d'études économiques* 78, 1462-1489 (10 2011) ; <https://doi.org/10.2307/41407068>.
- 684 J. Duggan, U. Sherman, R. Carbery, A. McDonnell, Gestion algorithmique et travail via application dans l'économie des petits boulots : un programme de recherche pour les relations de travail et la GRH. *Human Resource Management Journal* 30, 114–132 (2020) ; <https://doi.org/10.1111/1748-8583.12258>.
- 685 B. Bai, H. Dai, DJ Zhang, F. Zhang, H. Hu, Les impacts de l'attribution de tâches algorithmiques sur les perceptions d'équité et productivité : données probantes issues d'expériences sur le terrain. *Gestion des opérations de fabrication et de service : M & SOM* 24, 3060–3078 (2022) ; <https://doi.org/10.1287/msom.2022.1120>.
- 686 J. Howard, P. Schulte, Gestion des risques liés à l'IA sur le lieu de travail et avenir du travail. *American Journal of Industrial Medicine* 67, 980–993 (2024) ; <https://doi.org/10.1002/ajim.23653>.
- 687 A. Bernhardt, L. Kresge, R. Suleiman, Le lieu de travail axé sur les données et le plaidoyer en faveur des droits des travailleurs en matière de technologie. *Revue des relations industrielles et du travail* 76, 3–29 (2023) ; <https://doi.org/10.1177/00197939221131558>.

- 688 D. Acemoglu, P. Restrepo, Tâches, automatisation et montée des inégalités salariales aux États-Unis. *Econometrica : Journal of the Société économétrique* 90, 1973–2016 (2022) ; <https://doi.org/10.3982/ECTA19815>.
- 689 D. Acemoglu, Changement technique, inégalités et marché du travail. *Journal of Economic Literature* 40, 7–72 (2002) ; <https://doi.org/10.1257/0022051026976>.
- 690 DH Autor, Pourquoi y a-t-il encore autant d'emplois ? L'histoire et l'avenir de l'automatisation du lieu de travail. *The Journal of Economic Perspectives : A Journal of the American Economic Association* 29, 3–30 (2015) ; <https://doi.org/10.1257/jep.29.3.3>.
- 691 Ó. Afonso, R. Forte, Secteurs routiniers et non routiniers, automatisation des tâches et polarisation des salaires. *Appliquée Économie* (2023) ; <https://www.tandfonline.com/doi/abs/10.1080/00036846.2023.2280461>.
- 692 D. Acemoglu, J. Loebbing, « Automatisation et polarisation » (National Bureau of Economic Research, 2022) ; <https://doi.org/10.3386/w30528>.
- 693 D. Autor, « Appliquer l'IA pour reconstruire les emplois de la classe moyenne » (National Bureau of Economic Research, 2024) ; <https://doi.org/10.3386/w32140>.
- 694 L. Karabarbounis, Perspectives sur la part du travail. *The Journal of Economic Perspectives : A Journal of the American Economic Association* 38, 107–136 (2024) ; <https://doi.org/10.1257/jep.38.2.107>.
- 695 M. Ranaldi, Inégalités dans la composition du revenu. *Revue du revenu et de la richesse* 68, 139–160 (2022) ; <https://doi.org/10.1111/roiw.12503>.
- 696 T. Piketty, *Le Capital au XXIe siècle* (The Belknap Press de Harvard University Press, Cambridge Massachusetts, 2014) ; <https://www.hup.harvard.edu/books/9780674430006>.
- 697 B. Moll, L. Rachel, P. Restrepo, Croissance inégale : l'impact de l'automatisation sur les inégalités de revenus et de richesse. *Econometrica : Journal de la Société d'économétrie* 90, 2645–2683 (2022) ; <https://doi.org/10.3982/ECTA19417>.
- 698 C. Wang, M. Zheng, X. Bai, Y. Li, W. Shen, L'avenir des emplois en Chine sous l'impact de l'intelligence artificielle. *Lettres de recherche financière* 55, 103798 (2023) ; <https://doi.org/10.1016/j.frl.2023.103798>.
- 699 H. Firooz, Z. Liu, Y. Wang, « L'automatisation et l'essor des entreprises superstars » (Banque fédérale de San Francisco, 2022) ; <https://doi.org/10.24148/wp2022-05>.
- 700 CT Okolo, L'IA dans les pays du Sud : opportunités et défis vers une gouvernance plus inclusive, Brookings (2023) ; <https://www.brookings.edu/articles/ai-in-the-global-south-opportunities-and-challenges-towards-more-inclusive-governance/>.
- 701 A. Korinek, JE Stiglitz, « Intelligence artificielle, mondialisation et stratégies de développement économique » (National Bureau of Economic Research, 2021) ; <https://doi.org/10.3386/w28453>.
- 702 C. Alonso, A. Berg, S. Kothari, C. Papageorgiou, S. Rehman, « La révolution de l'IA provoquera-t-elle une grande divergence ? » (Fonds monétaire international, 2020) ; <https://www.imf.org/en/Publications/WP/Issues/2020/09/11/Will-the-AI-Revolution-Cause-a-Great-Divergence-49734>.
- 703 H. Nii-Aponsah, B. Verspagen, P. Mohnen, « La relocalisation induite par l'automatisation et ses implications potentielles pour les économies en développement » (UNU-MERIT, 2023) ; <https://ideas.repec.org/p/unm/unumer/2023018.html>.
- 704 J. Jacobs, « Comment l'IA générative change le secteur des services informatiques du Sud global » (Information Technology and Innovation Foundation, 2024) ; <https://itif.org/publications/2024/06/10/how-generative-ai-is-changing-the-global-souths-it-services-sector/>.
- 705 N. Otis, R. Clarke, S. Delecourt, D. Holtz, R. Koning, « L'impact inégal de l'IA générative sur la performance entrepreneuriale » (Harvard Business School, 2024) ; https://www.hbs.edu/ris/Publication%20Files/24-042_9ebd2f26-e292-404c-b858-3e883f0e11c0.pdf.
- 706 A. Merali, Lois d'échelle pour la productivité économique : preuves expérimentales en traduction assistée par LLM, arXiv [econ.GN] (2024) ; <http://arxiv.org/abs/2409.02391>.
- 707 K. McElheran, JF Li, E. Brynjolfsson, Z. Kroff, E. Dinlersoz, L. Foster, N. Zolas, AI Adoption en Amérique : qui, quoi, et Où. *Journal of Economics & Management Strategy* 33, 375–415 (2024) ; <https://doi.org/10.1111/jems.12576>.
- 708 K. Bonney, C. Breaux, C. Buffington, E. Dinlersoz, L. Foster, N. Goldschlag, J. Haltiwanger, Z. Kroff, K. Savage, L'impact de l'IA sur la main-d'œuvre : tâches versus emplois ? *Economics Letters* 244, 111971 (2024) ; <https://doi.org/10.1016/j.econlet.2024.111971>.
- 709 A. Kreacic, L. Uribe, J. Romeo, A. Lasater-Wille, R. Jesuthasan, S. Luong, « Comment l'IA générative transforme les entreprises et la société : le bon, le mauvais et tout le reste » (Oliver Wyman Forum, 2024) ; <https://www.oliverwymanforum.com/global-consumer-sentiment/how-will-ai-affect-global-economics.html>.
- 710 NG Otis, S. Delecourt, K. Cranney, R. Koning, « Données mondiales sur les écarts entre les sexes et l'IA générative » (Harvard Business School, 2024) ; <https://www.hbs.edu/faculty/Pages/item.aspx?num=66548>.

- 711* S. Jaffe, NP Shah, J. Butler, A. Farach, A. Cambon, B. Hecht, M. Schwarz, J. Teevan, « L'IA générative dans les lieux de travail réels » (Microsoft, 2024) ; <https://www.microsoft.com/en-us/research/publication/generative-ai-in-real-world-workplaces/> .
- 712* E. Wiles, L. Kraye, M. Abbadi, U. Awasthi, R. Kennedy, P. Mishkin, D. Sack, F. Candelon, GenAI comme exosquelette : preuves expérimentales sur les travailleurs du savoir utilisant GenAI sur de nouvelles compétences, Social Science Research Network (2024) ; <https://doi.org/10.2139/ssrn.4944588>.
- 713 A. Toner-Rodgers, Intelligence artificielle, découverte scientifique et innovation produit (2024) ; https://aidantr.github.io/files/AI_innovation.pdf.
- 714 T. Besiroglu, N. Emery-Xu, N. Thompson, Impacts économiques de la R&D augmentée par l'IA. Politique de recherche 53, 105037 (2024) ; <https://doi.org/10.1016/j.respol.2024.105037>.
- 715 S. McConnell, K. Fortson, D. Rotz, P. Schochet, P. Burkander, L. Rosenber, A. Matri, R. D'Amico, « Fournir des services publics de main-d'œuvre aux demandeurs d'emploi : résultats d'impact sur 15 mois des programmes pour adultes et travailleurs déplacés de la WIA » (Mathematica Policy Research, 2016) ; <https://mathematica.org/publications/providing-public-workforce-services-to-job-seekers-15-month-impact-findings-on-the-wia-adult>.
- 716 J. Furman, « Les politiques pour l'avenir du travail devraient être fondées sur son passé et son présent » (Economic Innovation Group, 2024) ; <https://eig.org/wp-content/uploads/2024/07/TAWP-Furman.pdf>.
- 717 A. Anthony, L. Sharma, E. Noor, « Promouvoir un programme plus global pour une intelligence artificielle digne de confiance » (Dotation Carnegie pour la paix internationale, 2024) ; <https://carnegieendowment.org/research/2024/04/advancing-a-more-global-agenda-for-trustworthy-artificial-intelligence?lang=fr> .
- 718 S. Ghosh, A. Caliskan, « ChatGPT perpétue les préjugés sexistes dans la traduction automatique et ignore les pronoms non genrés : résultats sur le bengali et cinq autres langues à faibles ressources » dans Actes de la conférence 2023 AAAI/ACM sur l'IA, l'éthique et la société (AIES '23) (Association for Computing Machinery, New York, NY, États-Unis, 2023), pp. 901–912 ; <https://doi.org/10.1145/3600211.3604672>.
- 719 C. Okorie, V. Marivate, Comment les experts africains en PNL relèvent les défis du droit d'auteur, de l'innovation et l'accès, Carnegie Endowment for International Peace (2024) ; <https://carnegieendowment.org/research/2024/04/how-african-nlp-experts-are-navigating-the-challenges-of-copyright-innovation-and-access?lang=fr> .
- 720 N. Maslej, L. Fattorini, E. Brynjolfsson, J. Etchemendy, K. Ligett, T. Lyons, J. Manyika, H. Ngo, JC Niebles, V. Parli, Y. Shoham, R. Wald, J. Clark, R. Perrault, « Rapport sur l'indice d'intelligence artificielle 2023 » (Comité directeur de l'indice d'intelligence artificielle, Institut pour l'IA centrée sur l'humain, Université de Stanford, 2023) ; <https://arxiv.org/pdf/2310.03715>.
- 721 N. Ahmed, M. Wahed, NC Thompson, L'influence croissante de l'industrie dans la recherche en IA. Science (New York, NY) 379, 884–886 (2023) ; <https://doi.org/10.1126/science.ade2420>.
- 722 S. Teleanu, J. Kurbalija, « Des voix numériques plus fortes en Afrique : construire une politique étrangère numérique africaine et Diplomatie » (Diplo, 2022) ; <https://www.diplomacy.edu/resource/report-stronger-digital-voices-from-africa/>.
- 723 T. Alsop, Expéditions estimées d'unités de traitement graphique (GPU) Nvidia H100 dans le monde en 2023, par Client, Statista (2024) ; <https://www.statista.com/statistics/1446564/nvidia-h100-gpu-shipments-by-customer/> .
- 724* Centres de données Google, Investir au Nebraska (2020) ; <https://www.google.com/intl/es/about/datacenters/locations/papillion/>.
- 725 Bureau du gouverneur Michael L. Parson, le gouverneur Parson annonce la sélection de Kansas City par Google pour un nouveau centre de données (2024) ; <https://governor.mo.gov/press-releases/archive/governor-parson-announces-googles-selection-kansas-city-new-data-center> .
- 726* Meta, « Centre de données de Meta à Prineville » (Meta, 2024) ; <https://datacenters.atmeta.com/wp-content/téléchargements/2024/10/Oregon-Prineville.pdf>.
- 727* Microsoft, Microsoft et le G42 annoncent une initiative globale d'écosystème numérique d'un milliard de dollars pour le Kenya, Stories (2024) ; <https://news.microsoft.com/2024/05/22/microsoft-and-g42-announce-1-billion-comprehensive-digital-ecosystem-initiative-for-kenya/> .
- 728 R. Zwetsloot, B. Zhang, N. Dreksler, L. Kahn, M. Anderjung, A. Dafeo, MC Horowitz, « Qualifié et mobile : enquête « Preuves des préférences des chercheurs en IA en matière d'immigration » dans les actes de la conférence 2021 de l'AAAI/ACM sur l'IA, l'éthique et la société (AIES '21) (Association for Computing Machinery, New York, NY, États-Unis, 2021), pp. 1050–1059 ; <https://doi.org/10.1145/3461702.3462617>.
- 729 meilleures universités, classement mondial des universités QS pour la science des données et l'intelligence artificielle 2024 (2024) ; <https://www.topuniversities.com/university-subject-rankings/data-science-artificial-intelligence>.

- 730 N. Maslej, L. Fattorini, R. Perrault, V. Parli, A. Reuel, E. Brynjolfsson, J. Etchemendy, K. Ligett, T. Lyons, J. Manyika, J. C. Niebles, Y. Shoham, R. Wald, J. Clark, « The AI Index 2024 Annual Report » (Institut pour l'IA centrée sur l'humain, Université de Stanford, 2024) ; <https://aiindex.stanford.edu/report/>.
- 731 N. Maslej, L. Fattorini, R. Perrault, V. Parli, A. Reuel, E. Brynjolfsson, J. Etchemendy, K. Ligett, T. Lyons, J. Manyika, J. C. Niebles, Y. Shoham, R. Wald, J. Clark, « The AI Index 2024 Annual Report » (Institut pour l'IA centrée sur l'humain, Université de Stanford, 2024) ; <https://aiindex.stanford.edu/report/>.
- 732 ML Gray, S. Suri, *Ghost Work : Comment empêcher la Silicon Valley de créer une nouvelle sous-classe mondiale* (Houghton Mifflin Harcourt, 2019) ; <https://ghostwork.info/>.
- 733 A. Arora, M. Barrett, E. Lee, E. Oborn, K. Prince, *Le risque et l'avenir de l'IA : biais algorithmiques, colonialisme des données et Marginalisation*. *Information et Organisation* 33 (2023) ; <https://doi.org/10.1016/j.infoandorg.2023.100478>.
- 734 CT Okolo, « Résoudre les inégalités mondiales dans le développement de l'IA » dans *Handbook of Critical Studies of Artificial Renseignement* (Edward Elgar Publishing, 2023), pp. 378–389 ; <https://www.elgaronline.com/edcollchap/book/9781803928562/book-part-9781803928562-40.xml>.
- 735 M. Miceli, T. Yang, A. Alvarado Garcia, J. Posada, SM Wang, M. Pohl, A. Hanna, *Documentation de la production de données Processus : une approche participative pour le travail sur les données*. *Actes de l'ACM sur l'interaction homme-machine* 6, 1–34 (2022) ; <https://doi.org/10.1145/3555623>.
- 736 D. Wang, S. Prabhat, N. Sambasivan, « De qui rêve l'IA ? À la recherche de l'aspiration dans l'annotation des données » dans *Conférence CHI sur les facteurs humains dans les systèmes informatiques (CHI '22)* (ACM, La Nouvelle-Orléans LA USA, 2022), pp. 1–16 ; <https://doi.org/10.1145/3491102.3502121>.
- 737 M. Steiger, TJ Bharucha, S. Venkatagiri, MJ Riedl, M. Lease, « Le bien-être psychologique des modérateurs de contenu : le travail émotionnel de la modération commerciale et les pistes pour améliorer le soutien » dans *Actes de la conférence CHI 2021 sur les facteurs humains dans les systèmes informatiques* (ACM, New York, NY, États-Unis, 2021) ; <https://doi.org/10.1145/3411764.3445092>.
- 738 MM AlEmadi, W. Zaghouani, « Bilan émotionnel et stratégies d'adaptation : analyser les effets de l'annotation des données sur les discours de haine » dans *Actes de l'atelier sur les questions juridiques et éthiques dans les technologies du langage humain @ LREC-COLING 2024* (2024), pp. 66–72 ; <https://aclanthology.org/2024.legal-1.10.pdf>.
- 739 S. Luccioni, Y. Jernite, E. Strubell, « Traitement gourmand en énergie : les watts déterminent le coût du déploiement de l'IA ? » dans la conférence 2024 de l'ACM sur l'équité, la responsabilité et la transparence (ACM, New York, NY, États-Unis, 2024) ; <https://doi.org/10.1145/3630106.3658542>.
- 740 B. Thormundsson, « Évolution de la concentration des talents liés à l'intelligence artificielle (IA) dans le monde de 2016 à 2023, par région » (Statista, 2024) ; <https://www.statista.com/statistics/1472183/ai-talent-concentration-change-percentage-by-region/>.
- 741 SV Bentley, CK Naughtin, MJ McGrath, JL Irons, PS Cooper, *La fracture numérique en action : comment les expériences de la technologie numérique façonnent les relations futures avec l'intelligence artificielle*. *AI and Ethics* 4, 901–915 (2024) ; <https://doi.org/10.1007/s43681-024-00452-3>.
- 742 Nigéria Ministère fédéral des communications, de l'innovation et de l'économie numérique, « Accélérer notre développement collectif « La prospérité grâce à l'efficacité technique : un plan stratégique pour le ministère fédéral des Communications, de l'Innovation et de l'Économie numérique » (2023) ; <https://fmcide.gov.ng/wp-content/uploads/2023/11/blueprint.pdf>.
- 743 Gouvernement américain, *Apportez vos compétences en IA aux États-Unis* (2023) ; <https://ai.gov/immigrate/>.
- 744 Gouvernement britannique, *Soutenir la prochaine génération de leaders de l'IA du monde entier* (2023) ; <https://www.great.gov.uk/campaign-site/ai-futures/>.
- 745 S. Pal, « D'où vient la main-d'œuvre européenne de l'IA ? Immigration, émigration et mouvement transfrontalier des talents de l'IA » (interface, 2024) ; <https://www.stiftung-nv.de/publications/where-is-europes-ai-workforce-coming-from>.
- 746 M. Mazumder, C. Banbury, X. Yao, B. Karlaş, WG Rojas, S. Damos, G. Damos, L. He, A. Parrish, HR Kirk, J. Quaye, C. Rastogi, D. Kiela, D. Jurado, D. Kanter, R. Mosquera, J. Ciro, ... VJ Reddi, « DataPerf : repères pour le développement de l'IA centrée sur les données » dans *37e Conférence internationale sur les systèmes de traitement de l'information neuronale (NeurIPS 2023)* (Curran Associates Inc., Red Hook, NY, États-Unis, 2024), pp. 5320–5347 ; <https://doi.org/10.5555/3666122.3666357>.
- 747 N. Guha, J. Nyarko, DE Ho, C. Ré, « Building GenAI Benchmarks: A Case Study in Legal Applications » dans *The Oxford Handbook on the Foundations and Regulation of Generative AI*, P. Hacker, A. Engel, S. Hammer, B. Mittelstadt, éd. (Oxford University Press, Oxford, Angleterre) ; https://neelguha.github.io/assets/pdf/building_genai_benchmarks_for_law_oxford_chapter.pdf.
- 748 E. Brynjolfsson, A. Ng, « La grande IA peut centraliser la prise de décision et le pouvoir, et c'est un problème » dans *Missing Links*

- dans *AI Governance*, B. Prud'homme, C. Régis, G. Farnadi, Eds. (UNESCO/MILA, 2023), pp. 65-87 ; <https://www.unesco.org/en/articles/missing-links-ai-governance>.
- 749 A. Korinek, J. Vipra, « Concentration de l'intelligence : mise à l'échelle et structure du marché en intelligence artificielle » (w33139, Bureau national de recherche économique, 2024) ; <https://doi.org/10.3386/w33139>.
- 750 Competition and Markets Authority, « Modèles de base de l'IA : rapport initial » (CMA, 2023) ; <https://www.gov.uk/government/publications/ai-foundation-models-initial-report>.
- 751 A. Chowdhery, S. Narang, J. Devlin, M. Bosma, G. Mishra, A. Roberts, P. Barham, H. W. Chung, C. Sutton, S. Gehrmann, P. Schuh, K. Shi, S. Tsvyashchenko, J. Maynez, A. Rao, P. Barnes, Y. Tay, ... N. Fiedel, PaLM : Mise à l'échelle de la modélisation du langage avec des chemins. *Journal of Machine Learning Research : JMLR* 24, 240 : 11324–240 : 11436 (2024).
- 752 X. Jin, D. Zhang, H. Zhu, W. Xiao, S.-W. Li, X. Wei, A. Arnold, X. Ren, « Préformation tout au long de la vie : adaptation continue « Language Models to Emerging Corpora » dans *Proceedings of BigScience Episode #5 – Atelier sur les défis et les perspectives de la création de grands modèles linguistiques* (Association for Computational Linguistics, Stroudsburg, PA, États-Unis, 2022), pp. 1–16 ; <https://doi.org/10.18653/v1/2022.bigscience-1.1>.
- 753 K. Gupta, B. Thérien, A. Ibrahim, ML Richter, QG Anthony, E. Belilovsky, I. Rish, T. Lesort, « Pré-entraînement continu de grands modèles de langage : comment réchauffer votre modèle ? » dans *Atelier sur les systèmes efficaces pour les modèles de base @ ICML2023* (2023) ; <https://openreview.net/pdf?id=pg7PUJe0TI>.
- 754 D. Luitse, Le pouvoir des plateformes dans l'IA : l'évolution des infrastructures cloud dans l'économie politique de l'intelligence artificielle. *Internet Policy Review* 13, 1–44 (2024) ; <https://doi.org/10.14763/2024.2.1768>.
- 755 C. Rikap, Variétés de systèmes d'innovation d'entreprise et leur interaction avec les systèmes mondiaux et nationaux : Stratégies d'Amazon, Facebook, Google et Microsoft pour produire et s'approprier l'intelligence artificielle. *Revue d'économie politique internationale*, 1–29 (2024) ; <https://doi.org/10.1080/09692290.2024.2365757>.
- 756 F. Richter, Amazon conserve son avance dans le cloud tandis que Microsoft s'en rapproche, *Statista* (2024) ; <https://www.statista.com/chart/18819/part-de-marche-mondiale-des-principaux-fournisseurs-de-services-d-infrastructure-cloud> .
- 757 P. Maham, S. Küspert, « Gouvernance de l'IA à usage général : une carte complète du manque de fiabilité, de l'utilisation abusive et Risques systémiques » (Stiftung Neue Verantwortung, 2023) ; <https://www.interface-eu.org/publications/governing-general-purpose-ai-comprehensive-map-unreliability-misuse-and-systemic-risks> .
- 758 G. Yu, G. Tan, H. Huang, Z. Zhang, P. Chen, R. Natella, Z. Zheng, A Survey on Failure Analysis and Fault Injection in AI Systèmes, *arXiv [cs.SE]* (2024) ; <http://arxiv.org/abs/2407.00125>.
- 759 F. Jimmy, Menaces émergentes : les derniers risques de cybersécurité et le rôle de l'intelligence artificielle dans l'amélioration des défenses de cybersécurité. *Revue internationale de recherche scientifique et de gestion* 9, 564–574 (2021) ; <https://doi.org/10.18535/ijrm/v9i2.ec01>.
- 760 Département du Trésor américain, Gestion des risques de cybersécurité spécifiques à l'intelligence artificielle dans le secteur financier Secteur des services. (2024) ; <https://home.treasury.gov/system/files/136/Managing-Artificial-Intelligence-Specific-Cybersecurity-Risks-In-The-Financial-Services-Sector.pdf> .
- 761 S. Trivedi, V. Aggarwal, R. Rastogi, « Améliorer la puissance des systèmes cyberphysiques activés par l'IA » dans *Solutions d'intelligence artificielle pour les systèmes cyberphysiques* (Auerbach Publications, Boca Raton, éd. 1, 2024), pp. 1–39 ; <https://doi.org/10.1201/9781032694375-1>.
- 762 ID Raji, S. Costanza-Chock, J. Buolamwini, « Le changement venu de l'extérieur : vers des audits crédibles des systèmes d'IA par des tiers », dans *Les chaînons manquants dans la gouvernance de l'IA*, B. Prud'homme, C. Régis, G. Farnadi, éd. (UNESCO/MILA, 2023), pp. 4–26 ; <https://www.unesco.org/fr/articles/les-chainons-manquants-dans-la-gouvernance-de-l-ia>.
- 763 M. Stein, M. Gandhi, T. Kriecherbauer, A. Oueslati, R. Trager, « Organismes publics ou privés : qui doit réaliser des évaluations et des audits d'IA avancés ? Une logique en trois étapes basée sur des études de cas d'industries à haut risque » (Oxford Martin AI Governance Initiative, 2024) ; <https://www.oxfordmartin.ox.ac.uk/publications/public-vs-private-bodies-who-should-run-advanced-ai-evaluations-and-audits-a-three-step-logic-based-on-case-studies-of-high-risk-industries> .
- 764 AJ Grotto, J. Dempsey, « Divulgarion et gestion des vulnérabilités pour les systèmes d'IA/ML : un document de travail avec des recommandations politiques » (Stanford Geopolitics, Technology, and Governance Cyber Policy Center, 2021) ; <https://doi.org/10.2139/ssrn.3964084>.
- 765 Y. Hong, J. Lian, L. Xu, J. Min, Y. Wang, LJ Freeman, X. Deng, Perspectives statistiques sur la fiabilité des Systèmes de renseignement. *Ingénierie de la qualité* 35, 56–78 (2023) ; <https://doi.org/10.1080/08982112.2022.2089854>.
- 766 T. Aguirre, Sur les laboratoires et les usines : cartographie de la manière dont les alliances, les acquisitions et les lois antitrust façonnent l'industrie de l'IA de pointe, *arXiv [econ.GN]* (2024) ; <http://arxiv.org/abs/2406.01722>.

- 767 B. Martens, « Pourquoi l'intelligence artificielle crée des défis fondamentaux pour la politique de la concurrence » (16/2024, Note Bruegel, 2024) ; <https://hdl.handle.net/10419/302296>.
- 768 Agence américaine de protection de l'environnement, « Calculateur d'équivalences de gaz à effet de serre - Calculs et références » (EPA, 2024) ; <https://www.epa.gov/energy/greenhouse-gas-equivalencies-calculator-calculations-and-references> .
- 769* Gemma Team, M. Rivière, S. Pathak, PG Sessa, C. Hardin, S. Bhupatiraju, L. Hussenot, T. Mesnard, B. Shahriari, A. Ramé, J. Ferret, P. Liu, P. Tafti, A. Friesen, M. Casbon, S. Ramos, R. Kumar, ... A. Andreev, Gemma 2 : Améliorer les modèles de langage ouverts à une taille pratique, arXiv [cs.CL] (2024) ; <http://arxiv.org/abs/2408.00118>.
- 770 D. Donnellan, A. Lawrence, D. Bizo, P. Judge, J. O'Brien, J. Davis, M. Smolaks, J. Williams-George, R. Weinschenk, « Enquête mondiale sur les centres de données 2024 de l'Uptime Institute » (Uptime Institute, 2024) ; <https://uptimeinstitute.com/resources/research-and-reports/uptime-institute-global-data-center-survey-results-2024> .
- 771 V. Rozite, E. Bertoli, B. Reidenbach, « Centres de données et réseaux de transmission de données » (Agence internationale de l'énergie, 2023) ; <https://www.iea.org/energy-system/buildings/data-centres-and-data-transmission-networks>.
- 772 L. Burdette, P. Brodsky, P. Christian, J. Hjembo, A. Mauldin, T. Stronge, M. Tan, J. Velandia, « The State of the Réseau Édition 2023 » (TeleGeography, 2023) ; <https://www2.telegeography.com/hubfs/LP-Assets/Ebooks/state-of-the-network-2023.pdf> .
- 773 R. Schwartz, J. Dodge, N. A. Smith, O. Etzioni, G. A. Green. Communications de l'ACM 63, 54–63 (2020); <https://doi.org/10.1145/3381831>.
- 774 LH Kaack, PL Donti, E. Strubell, G. Kamiya, F. Creutzig, D. Rolnick, Aligner l'intelligence artificielle avec l'atténuation du changement climatique. Nature Climate Change 12, 518–527 (2022) ; <https://doi.org/10.1038/s41558-022-01377-7>.
- 775 E. Zelikman, Y. Wu, J. Mu, N. Goodman, « STaR : Bootstrapping Reasoning With Reasoning » dans Advances in Neural Information Processing Systems (NeurIPS 2022) (La Nouvelle-Orléans, LA, États-Unis, 2022) vol. 35, pp. 15476–15488 ; https://proceedings.neurips.cc/paper_files/paper/2022/file/639a9a172c044fbb64175b5fad42e9a5-Paper-Conference.pdf .
- 776* T. Wu, J. Lan, W. Yuan, J. Jiao, J. Weston, S. Sukhbaatar, Thinking LLMs : Instruction générale avec Génération de pensée, arXiv [cs.CL] (2024) ; <http://arxiv.org/abs/2410.10630>.
- 777 L. Long, R. Wang, R. Xiao, J. Zhao, X. Ding, G. Chen, H. Wang, « Sur la génération de données synthétiques pilotée par les LLM, Conservation et évaluation : une enquête » dans Findings of the Association for Computational Linguistics ACL 2024 (Association for Computational Linguistics, Stroudsburg, PA, États-Unis, 2024), pp. 11065–11082 ; <https://doi.org/10.18653/v1/2024.findings-acl.658>.
- 778 N. Alder, K. Ebert, R. Herbrich, P. Hacker, IA, climat et transparence : opérationnaliser et améliorer l'IA Acte, arXiv [cs.CY] (2024) ; <http://arxiv.org/abs/2409.07471>.
- 779* AS Luccioni, A. Hernandez-Garcia, Comptage du carbone : une étude des facteurs influençant les émissions des machines Apprentissage, arXiv [cs.LG] (2023) ; <http://arxiv.org/abs/2302.08476>.
- 780* Google, « Rapport environnemental 2024 » (2024) ; <https://www.gstatic.com/gumdrop/sustainability/google-2024-rapport-environnemental.pdf>.
- 781 Baidu, « Rapport environnemental, social et de gouvernance 2023 de Baidu » (2023) ; <https://esg.baidu.com/Uploads/File/2024/05/17/Baidu%202023%20Environmental,%20Social%20and%20Governance%20Report.20240517150706.pdf>.
- 782 EPRI, « Powering Intelligence : Analyse de l'intelligence artificielle et de la consommation énergétique des centres de données » (2024) ; <https://www.epri.com/research/products/00000003002028905>.
- 783 G. Guidi, F. Dominici, J. Gilmour, K. Butler, E. Bell, S. Delaney, FJ Bargagli-Stoffi, Charge environnementale des centres de données des États-Unis à l'ère de l'intelligence artificielle, arXiv [cs.CY] (2024) ; <http://arxiv.org/abs/2411.09786>.
- 784 Agence internationale de l'énergie, « World Energy Outlook 2024 » (AIE, 2024) ; <https://www.iea.org/reports/world-energy-outlook-2024> .
- 785 Office central des statistiques d'Irlande, « Consommation d'électricité mesurée dans les centres de données 2023 » (CSO, 2024) ; <https://www.cso.ie/en/releasesandpublications/ep/p-dcmecl/datacentresmeteredelectricityconsumption2023/>.
- 786 PGIM Real Estate, « Extrait de Global Data Centers Americas » (2021) ; https://cdn.pfcdn.com/cms/pgim-real-estate/sites/default/files/2021-01/Global%20Data%20Centers-US_February%202021_PGIM.pdf .
- 787 Bureau de la politique du ministère de l'Énergie des États-Unis, « Ressources énergétiques propres pour répondre à la demande en électricité des centres de données » (DOE, 2024) ; <https://www.energy.gov/policy/articles/clean-energy-resources-meet-data-center-electricity-demand> .

- 788* Constellation, Constellation va lancer le Crane Clean Energy Center, rétablissant les emplois et alimentant le réseau en électricité sans carbone (2024) ; <https://www.constellationenergy.com/newsroom/2024/Constellation-to-Launch-Crane-Clean-Energy-Center-Restoring-Jobs-and-Carbon-Free-Power-to-The-Grid.html> .
- 789 Talen Energy Corporation, « Libérer la valeur » (2024) ; <https://ir.talenenergy.com/static-files/f02c44a9-d2dc-45c1-9331-eee1495f7d2d>.
- 790 Commission fédérale de réglementation de l'énergie des États-Unis, ordonnance rejetant les modifications apportées à l'accord de service d'interconnexion. FERC (2024) ; https://elibrary.ferc.gov/eLibrary/filelist?accession_number=20241101-3061&optimisé=faux.
- 791* M. Terrell, Nouvel accord sur l'énergie nucléaire propre avec Kairos Power, Google (2024) ; <https://blog.google/outreach-initiatives/sustainability/google-kairos-power-nuclear-energy-agreement/>.
- 792 LM Krall, AM Macfarlane, RC Ewing, Déchets nucléaires provenant de petits réacteurs modulaires. Actes de la conférence nationale Académie des sciences des États-Unis d'Amérique 119, e2111833119 (2022) ; <https://doi.org/10.1073/pnas.2111833119>.
- 793 J. Dodge, T. Prewitt, R. Tachet des Combes, E. Odmark, R. Schwartz, E. Strubell, AS Luccioni, NA Smith, N. DeCario, W. Buchanan, « Mesurer l'intensité carbone de l'IA dans les instances Cloud » dans la conférence 2022 de l'ACM sur l'équité, la responsabilité et la transparence (ACM, New York, NY, États-Unis, 2022) ; <https://doi.org/10.1145/3531146.3533234>.
- 794 P. Hacker, Régulation durable de l'IA, arXiv [cs.CY] (2023) ; <http://arxiv.org/abs/2306.00292>.
- 795* Meta, « Rapport de développement durable 2024 » (2024) ; <https://sustainability.atmeta.com/wp-content/téléchargements/2024/08/Meta-2024-Rapport-Durabilité.pdf>.
- 796* Amazon, « Rapport sur le développement durable d'Amazon » (2024) ; <https://sustainability.aboutamazon.com/2023-amazon-rapport-de-durabilite.pdf>.
- 797 AN Achanta, P. Erickson, E. Haites, M. Lazarus, N. Pandey, N. Pahuja, S. Seres, R. Spalding-Fecher, R. Tewari, « Évaluation de l'impact du mécanisme de développement propre » (Groupe de haut niveau sur le dialogue politique sur le MDP, 2012) ; https://www.cdmpolicydialogue.org/research/1030_impact.pdf.
- 798 J. Rasley, S. Rajbhandari, O. Ruwase, Y. He, « DeepSpeed : les optimisations du système permettent la formation à l'apprentissage profond Modèles avec plus de 100 milliards de paramètres » dans les actes de la 26e conférence internationale ACM SIGKDD sur la découverte des connaissances et l'exploration des données (ACM, New York, NY, États-Unis, 2020) ; <https://doi.org/10.1145/3394486.3406703>.
- 799 W. Kwon, Z. Li, S. Zhuang, Y. Sheng, L. Zheng, CH Yu, J. Gonzalez, H. Zhang, I. Stoica, « Efficient Memory « Gestion des modèles de langage volumineux avec PagedAttention » dans Actes du 29e Symposium sur les principes des systèmes d'exploitation (ACM, New York, NY, États-Unis, 2023), pp. 611–626 ; <https://doi.org/10.1145/3600006.3613165>.
- 800 HD Saunders, Le postulat de Khazzoom-Brookes et la croissance néoclassique. *The Energy Journal* 13, 131–148 (1992) ; <http://www.jstor.org/stable/41322471>.
- 801 G. Kamiya, VC Coroamă, « Consommation d'énergie des centres de données – Un examen critique » (IEA 4E TCP Electronic Devices and Annexe des réseaux (EDNA)).
- 802 Agence internationale de l'énergie, « Suivi des progrès en matière d'énergie propre 2023 » (AIE, 2023) ; <https://www.iea.org/reports/tracking-clean-energy-progress-2023>.
- 803 E. Halper, Dans un contexte de demande explosive, l'Amérique est à court d'énergie, *Washington Post* (2024) ; <https://www.washingtonpost.com/business/2024/03/07/ai-data-centers-power/>.
- 804 Commission européenne, Centre commun de recherche, G. Kamiya, P. Bertoldi, Consommation d'énergie dans les centres de données et Réseaux de communication à large bande dans l'UE (Office des publications de l'Union européenne, 2024) ; <https://doi.org/10.2760/706491>.
- 805 J. Koomey, E. Masanet, Does Not Compute : éviter les pièges de l'évaluation des impacts énergétiques et carbone d'Internet. *Joule* 5, 1625–1628 (2021) ; <https://doi.org/10.1016/j.joule.2021.05.007>.
- 806 E. Masanet, A. Shehabi, N. Lei, S. Smith, J. Koomey, Réévaluation des estimations de la consommation d'énergie des centres de données mondiaux. *Science* (New York, NY) 367, 984–986 (2020) ; <https://doi.org/10.1126/science.aba3758>.
- 807 D. Rolnick, PL Donti, LH Kaack, K. Kochanski, A. Lacoste, K. Sankaran, AS Ross, N. Milojevic-Dupont, N. Jaques, A. Waldman-Brown, AS Luccioni, T. Maharaj, ED Sherwin, SK Mukkavilli, KP Kording, CP Gomes, AY Ng,... Y. Bengio, Lutter contre le changement climatique grâce à l'apprentissage automatique. *ACM Computing Surveys* 55, 1–96 (2023) ; <https://doi.org/10.1145/3485128>.
- 808 U. Gupta, YG Kim, S. Lee, J. Tse, H.-HS Lee, G.-Y. Wei, D. Brooks, C.-J. Wu, À la poursuite du carbone : l'insaisissable Empreinte environnementale de l'informatique. *IEEE Micro* 42, 37–47 (2022) ; <https://doi.org/10.1109/mm.2022.3163226>.
- 809* Intel, « Rapport sur la responsabilité d'entreprise 2023-24 » (2024) ;

- <https://csrreportbuilder.intel.com/pdfbuilder/pdfs/CSR-2023-24-Full-Report.pdf>.
- 810 Agence européenne pour l'environnement, « Utilisation de l'eau en Europe : la quantité et la qualité sont confrontées à de grands défis » (AEE, 2018) ; <https://www.eea.europa.eu/signals-archived/signals-2018-content-list/articles/water-use-in-europe-2014>.
- 811 Taiwan Semiconductor Manufacturing Company, « Rapport sur le développement durable TSMC 2023 » (TSMC, 2024) ; https://esg.tsmc.com/en-US/file/public/e-all_2023.pdf.
- 812 P. Li, J. Yang, MA Islam, S. Ren, Rendre l'IA moins « assoiffée » : découvrir et traiter l'empreinte hydrique secrète des modèles d'IA, arXiv [cs.LG] (2023) ; <http://arxiv.org/abs/2304.03271>.
- 813 Nations Unies, Le droit humain à l'eau et à l'assainissement : Résolution A/RES/64/292 adoptée par l'Assemblée générale le 28 juillet 2010 (2010) ; <https://documents.un.org/doc/undoc/gen/n09/479/35/pdf/n0947935.pdf>.
- 814 Le Parlement européen et le Conseil de l'Union européenne, Directive (UE) 2023/1791 du Parlement européen et du Conseil Parlement européen et du Conseil sur l'efficacité énergétique et modifiant le règlement (UE) 2023/955 (refonte) (texte présentant de l'intérêt pour l'EEE). (2023) ; <https://eur-lex.europa.eu/eli/dir/2023/1791/oj>.
- 815 Y. Jin, P. Behrens, A. Tukker, L. Scherer, Utilisation de l'eau dans les technologies électriques : une méta-analyse mondiale. Revue des énergies renouvelables et durables 115, 109391 (2019) ; <https://doi.org/10.1016/j.rser.2019.109391>.
- 816 H. Zhai, ES Rubin, EJ Grol, AC O'Connell, Z. Wu, EG Lewis, Modernisation du refroidissement à sec dans les centrales électriques à combustible fossile existantes dans une région soumise à un stress hydrique : compromis en matière d'économies d'eau, de coûts et de déficits de capacité. Applied Energy 306, 117997 (2022) ; <https://doi.org/10.1016/j.apenergy.2021.117997>.
- 817 VG Gude, Consommation et récupération d'énergie dans l'osmose inverse. Dessalement et traitement de l'eau 36, 239–260 (2011) ; <https://doi.org/10.5004/dwt.2011.2534>.
- 818 Ministère australien de l'environnement et de l'énergie, « Fiche d'information sur le CVC : Co et trigénération » (2013) ; <https://www.energy.gov.au/sites/default/files/hvac-factsheet-co-tri-generation.pdf>.
- 819 Office of Fossil Energy, « Stratégie de l'hydrogène : permettre une économie à faibles émissions de carbone » (Département américain de l'énergie, 2020) ; https://www.energy.gov/sites/prod/files/2020/07/f76/USDOE_FE_Hydrogen_Strategy_July2020.pdf.
- 820 H. Nissenbaum, La vie privée dans son contexte : technologie, politique et intégrité de la vie sociale (Stanford University Press, Palo Alto, Californie, 2009) ; <http://www.sup.org/books/title/?id=8862>.
- 821 L. Bourtole, V. Chandrasekaran, CA Choquette-Choo, H. Jia, A. Travers, B. Zhang, D. Lie, N. Papernot, « Machine « Désapprendre » dans le Symposium IEEE 2021 sur la sécurité et la confidentialité (SP) (IEEE, virtuel, 2021), pp. 141–159 ; <https://doi.org/10.1109/SP40001.2021.00019>.
- 822 Organisation de coopération et de développement économiques, « IA, gouvernance des données et confidentialité » (OCDE, 2024) ; <https://doi.org/10.1787/2476b1a4-en>.
- 823 Comité européen de la protection des données, « Rapport sur les travaux entrepris par le groupe de travail ChatGPT » (EDPB, 2024) ; https://www.edpb.europa.eu/our-work-tools/our-documents/other/report-work-undertaken-chatgpt-taskforce_en.
- 824 DJ Solove, Intelligence artificielle et confidentialité. Florida Law Review (à paraître en janvier 2025) ; https://papers.ssrn.com/sol3/papers.cfm?abstract_id=4713111.
- 825 Parlement du Royaume-Uni, Loi sur la protection des données de 2018, article 46 : Droit de rectification. (2018) ; <https://www.legislation.gov.uk/ukpga/2018/12/section/46>.
- 826 Groupe de travail sur la coopération internationale en matière d'application de la loi du GPA, « Déclaration conjointe sur le grattage de données et la protection de la vie privée » (Information Commissioner's Office, 2023) ; <https://ico.org.uk/media/about-the-ico/documents/4026232/joint-statement-data-scraping-202308.pdf>.
- 827 N. Carlini, D. Ippolito, M. Jagielski, K. Lee, F. Tramer, C. Zhang, « Quantification de la mémorisation dans les modèles de langage neuronal » dans 11e Conférence internationale sur les représentations de l'apprentissage (ICLR 2023) (Kigali, Rwanda, 2022) ; https://openreview.net/forum?id=TatRHT_1cK.
- 828 Y. Chen, E. Mendes, S. Das, W. Xu, A. Ritter, Les modèles linguistiques peuvent-ils être chargés de protéger les informations personnelles ?, arXiv [cs.CL] (2023) ; <http://arxiv.org/abs/2310.02224>.
- 829 R. Shokri, M. Stronati, C. Song, V. Shmatikov, « Attaques par inférence d'appartenance contre les modèles d'apprentissage automatique » dans Symposium IEEE 2017 sur la sécurité et la confidentialité (SP) (IEEE, San Jose, CA, États-Unis, mai 2017), pp. 3–18 ; <https://doi.org/10.1109/SP.2017.41>.
- 830 M. Fredrikson, S. Jha, T. Ristenpart, « Attaques par inversion de modèle qui exploitent les informations de confiance et les Contre-mesures » dans les actes de la 22e conférence ACM SIGSAC sur la sécurité informatique et des communications (CCS '15) (Association for Computing Machinery, New York, NY, États-Unis, 2015), pp. 1322–1333 ; <https://doi.org/10.1145/2810103.2813677>.
- 831 M. Duan, A. Suri, N. Mireshghallah, S. Min, W. Shi, L. Zettlemoyer, Y. Tsvetkov, Y. Choi, D. Evans, H. Hajishirzi, Do

- Les attaques par inférence d'appartenance fonctionnent-elles sur de grands modèles de langage ?, arXiv [cs.CL] (2024) ; <http://arxiv.org/abs/2402.07841>.
- 832 N. Carlini, F. Tramèr, E. Wallace, M. Jagielski, A. Herbert-Voss, K. Lee, A. Roberts, T. Brown, D. Song, Ú. Erlingsson, A. Oprea, C. Raffel, « Extraction de données de formation à partir de grands modèles de langage » dans 30e symposium sur la sécurité USENIX (USENIX Security 21) (USENIX Association, 2021), pp. 2633–2650 ; <https://www.usenix.org/conference/usenixsecurity21/presentation/carlini-extracting>.
- 833 N. Carlini, J. Hayes, M. Nasr, M. Jagielski, V. Sehwag, F. Tramèr, B. Balle, D. Ippolito, E. Wallace, « Extraction de données d'entraînement à partir de modèles de diffusion » dans 32e Symposium sur la sécurité USENIX (USENIX Security 23) (USENIX Association, Anaheim, CA, 2023), pp. 5253–5270 ; <https://www.usenix.org/conference/usenixsecurity23/presentation/carlini>.
- 834 W. Shi, A. Ajith, M. Xia, Y. Huang, D. Liu, T. Blevins, D. Chen, L. Zettlemoyer, « Détection de données de pré-entraînement à partir de grands modèles linguistiques » dans la 12e Conférence internationale sur les représentations d'apprentissage (ICLR 2024) (Vienne, Autriche, 2023) ; <https://openreview.net/forum?id=zWqr3MQUNs>.
- 835 N. Lukas, A. Salem, R. Sim, S. Tople, L. Wutschitz, S. Zanella-Béguelin, « Analyse des fuites d'informations personnelles identifiables dans les modèles linguistiques » dans Symposium IEEE 2023 sur la sécurité et la confidentialité (SP) (IEEE, 2023), pp. 346–363 ; <https://doi.org/10.1109/SP46215.2023.10179300>.
- 836 S. Longpre, R. Mahari, AN Lee, CS Lund, H. Oderinwale, W. Brannon, N. Saxena, N. Obeng-Marnu, T. South, CJ Hunter, K. Klyman, C. Klamm, H. Schoelkopf, N. Singh, M. Cherep, AM Anis, A. Dinh, ... A. Pentland, « Consentement en crise : le déclin rapide des données communes de l'IA » dans 38e Conférence sur les ensembles de données et les repères des systèmes de traitement de l'information neuronale (2024) ; <https://openreview.net/pdf?id=66PcEzKf95>.
- 837* K. Saab, T. Tu, W.-H. Weng, R. Tanno, D. Stutz, E. Wulczyn, F. Zhang, T. Strother, C. Park, E. Vedadi, JZ Chaves, S.-Y. Hu, M. Schaekermann, A. Kamath, Y. Cheng, DGT Barrett, C. Cheung, ... V. Natarajan, « Capacités des modèles Gemini en médecine » (Google Deepmind, 2024) ; <http://arxiv.org/abs/2404.18416>.
- 838 P. Lewis, E. Perez, A. Piktus, F. Petroni, V. Karpukhin, N. Goyal, H. Küttler, M. Lewis, W.-T. Yih, T. Rocktäschel, S. Riedel, D. Kiela, « Génération augmentée par récupération pour les tâches de traitement du langage naturel à forte intensité de connaissances » dans 34e conférence sur les systèmes de traitement de l'information neuronale (NeurIPS 2020) (Curran Associates, Inc., Vancouver, Canada, 2020) vol. 33, pp. 9459–9474 ; <https://proceedings.neurips.cc/paper/2020/hash/6b493230205f780e1bc26945df7481e5-Résumé.html>.
- 839 V. Karpukhin, B. Oguz, S. Min, P. Lewis, L. Wu, S. Edunov, D. Chen, W.-T. Yih, « Récupération de passages denses pour les « Réponses aux questions de domaine » dans les actes de la conférence 2020 sur les méthodes empiriques en traitement du langage naturel (EMNLP) (Association for Computational Linguistics, Stroudsburg, PA, États-Unis, 2020), pp. 6769–6781 ; <https://doi.org/10.18653/v1/2020.emnlp-main.550>.
- 840 O. Ram, Y. Levine, I. Dalmedigos, D. Muhlgay, A. Shashua, K. Leyton-Brown, Y. Shoham, Modèles linguistiques augmentés par recherche en contexte. Transactions of the Association for Computational Linguistics 11, 1316–1331 (2023) ; https://doi.org/10.1162/tacl_a_00605.
- 841* T. Gunter, Z. Wang, C. Wang, R. Pang, A. Narayanan, A. Zhang, B. Zhang, C. Chen, C.-C. Chiu, D. Qiu, D. Gopinath, D. A. Yap, D. Yin, F. Nan, F. Weers, G. Yin, H. Huang, ... Z. Ren, Modèles de langage de la Fondation Apple Intelligence, arXiv [cs.AI] (2024) ; <http://arxiv.org/abs/2407.21075>.
- 842 S. Arora, P. Lewis, A. Fan, J. Kahn, C. Ré, Raisonnement sur les données publiques et privées dans les systèmes basés sur la recherche. Transactions de l'Association pour la linguistique computationnelle 11, 902–921 (2023) ; https://doi.org/10.1162/tacl_a_00580.
- 843 G. Zyskind, T. South, A. Pentland, « N'oubliez pas la récupération privée : recherche de similarité privée distribuée pour les grands modèles de langage » dans Actes du cinquième atelier sur la confidentialité dans le traitement du langage naturel (2024), pp. 7–19 ; <https://aclanthology.org/2024.privatenlp-1.2.pdf>.
- 844 Centre national de cybersécurité du Royaume-Uni, Agence de cybersécurité et de sécurité des infrastructures des États-Unis, Agence de sécurité nationale, Federal Bureau of Investigation, Centre australien de cybersécurité de la Direction des signaux de l'Australie, Centre canadien pour la cybersécurité, Centre national de cybersécurité de la Nouvelle-Zélande, CSIRT du gouvernement chilien, Agence nationale de cybersécurité et de sécurité de l'information de la République tchèque, Autorité des systèmes d'information d'Estonie, Centre national de cybersécurité d'Estonie, Agence française de cybersécurité, Office fédéral allemand de la sécurité de l'information, Direction nationale israélienne de la cybersécurité, Agence nationale italienne de cybersécurité, Centre national japonais de préparation aux incidents et de stratégie pour la cybersécurité, Secrétariat japonais de la politique scientifique, technologique et d'innovation, Cabinet Office, ... Agence de cybersécurité de Singapour, « Lignes directrices pour le développement de systèmes d'IA sécurisés » (gouvernement du Royaume-Uni, 2023) ; <https://www.ncsc.gov.uk/files/Guidelines-for-secure-AI-system-development.pdf>.
- 845 M. Kosinski, D. Stillwell, T. Graepel, Les traits et attributs privés sont prévisibles à partir des enregistrements numériques du comportement humain. Actes de l'Académie nationale des sciences des États-Unis d'Amérique 110, 5802–5805

- (2013) ; <https://doi.org/10.1073/pnas.1218772110>.
- 846 R. Staab, M. Vero, M. Balunovic, M. Vechev, « Au-delà de la mémorisation : violation de la vie privée via l'inférence avec de grands modèles linguistiques » dans la 12e Conférence internationale sur les représentations d'apprentissage (ICLR 2024) (Vienne, Autriche, 2023) ; <https://openreview.net/forum?id=kmn0BhQk7p>.
- 847 N. Miresghallah, M. Antoniak, Y. More, Y. Choi, G. Farnadi, « Ne faites confiance à aucun bot : découverte de divulgations personnelles dans les conversations Human-LLM dans la nature » dans Première conférence sur la modélisation du langage (2024) ; <https://openreview.net/pdf?id=tlpWtMYkzU>.
- 848* J. Lamb, G. Israelstam, R. Agarwal, S. Bhasker, « L'IA générative dans le secteur de la santé : tendances d'adoption et prochaines étapes » (McKinsey & Company, 2024) ; <https://www.mckinsey.com/industries/healthcare/our-insights/generative-ai-in-healthcare-adoption-trends-and-whats-next>.
- 849 La Commission fédérale du commerce et le rapport du personnel de la FTC révèlent que les grandes entreprises de médias sociaux et de streaming vidéo ont engagé dans une surveillance à grande échelle des utilisateurs avec des contrôles de confidentialité laxistes et des garanties inadéquates pour les enfants et les adolescents (2024) ; <https://www.ftc.gov/news-events/news/press-releases/2024/09/ftc-staff-report-finds-large-social-media-video-streaming-companies-have-engaged-vast-surveillance>.
- 850 Federal Trade Commission, la FTC déclare que les employés de Ring ont surveillé illégalement les clients et n'ont pas réussi à empêcher les pirates de prendre le contrôle des caméras des utilisateurs (2023) ; <https://www.ftc.gov/news-events/news/press-releases/2023/05/ftc-says-ring-employees-illegally-surveilled-customers-failed-stop-hackers-taking-control-users>.
- 851* J. Ho, W. Chan, C. Saharia, J. Whang, R. Gao, A. Gritsenko, DP Kingma, B. Poole, M. Norouzi, DJ Fleet, T. Salimans, Imagen Video : Génération de vidéos haute définition avec des modèles de diffusion, arXiv [cs.CV] (2022) ; <http://arxiv.org/abs/2210.02303>.
- 852* Équipe Reka, A. Ormazabal, C. Zheng, C. de M. d'Autume, D. Yogatama, D. Fu, D. Ong, E. Chen, E. Lamprecht, H. Pham, I. Ong, K. Aleksiev, L. Li, M. Henderson, M. Bain, M. Artetxe, N. Relan, ... Z. Xie, Reka Core, Flash et Edge : une série de puissants modèles de langage multimodaux, arXiv [cs.CL] (2024) ; <http://arxiv.org/abs/2404.12387>.
- 853 S. Adler, Z. Hitzig, S. Jain, C. Brewer, W. Chang, R. DiResta, E. Lazzarin, S. McGregor, W. Seltzer, D. Siddarth, N. Soliman, T. South, C. Spelliscy, M. Sporny, V. Srivastava, J. Bailey, B. Christian, ... T. Zick, Personhood Credentials: Artificial Intelligence and the Value of Privacy-Preserving Tools to Distinguish Who Is Real Online, arXiv [cs.CY] (2024) ; <http://arxiv.org/abs/2408.07892>.
- 854 B. Auxier, L. Rainie, M. Anderson, A. Perrin, M. Kumar, E. Turner, « Les Américains et la vie privée : préoccupés, confus et ressentant un manque de contrôle sur leurs informations personnelles » (Pew Research Center, 2019) ; <https://www.pewresearch.org/internet/2019/11/15/americans-and-privacy-concerned-confused-and-feeling-lack-of-control-over-their-personal-information/>.
- 855* IBM, « Coût d'une violation de données 2024 » (2024) ; <https://www.ibm.com/reports/data-breach>.
- 856 S. Min, S. Gururangan, E. Wallace, W. Shi, H. Hajishirzi, NA Smith, L. Zettlemoyer, « Modèles de langage SILO : isoler le risque juridique dans un magasin de données non paramétrique » dans NeurIPS 2023 Workshop on Distribution Shifts (DistShift) (La Nouvelle-Orléans, LA, États-Unis, 2023) ; <https://openreview.net/forum?id=z03bW0doni>.
- 857 US Copyright Office, « Copyright et intelligence artificielle » (2024) ; <https://www.copyright.gov/ai/>.
- 858 P. Burger, La Convention de Berne : son histoire et son rôle clé dans l'avenir. *Journal of Law and Technology* 3, 1–70 (1988) ; <https://heinonline.org/HOL/P?h=hein.journals/jlawtecy3&i=9>.
- 859 LR Patterson, C. Joyce, Le droit d'auteur en 1791 : essai sur le point de vue des fondateurs sur le pouvoir du droit d'auteur accordé au Congrès dans l'article I, section 8, clause 8 de la Constitution américaine. *Emory Law Journal* (2003) ; https://heinonline.org/hol-cgi-bin/get_pdf.cgi?handle=hein.journals/emlj52§ion=25.
- 860 Le Bureau du Conseil de révision des lois de la Chambre des représentants des États-Unis, « Limitations des droits exclusifs Droits : Utilisation équitable. Sec. 107 » dans le Code des États-Unis, édition 2006, supplément 4, titre 17 - Droits d'auteur (US Government Publishing Office, éd. 2010, 2010) ; <https://www.govinfo.gov/app/details/USCODE-2010-titre17/USCODE-2010-titre17-chap1-sec107>.
- 861 Parlement européen, Direction générale des politiques internes de l'Union, E. Rosati, L'exception pour l'exploration de textes et de données (TDM) dans la proposition de directive sur le droit d'auteur dans le marché unique numérique – Aspects techniques (Parlement européen, 2018).
- 862 Système de base de données de traduction juridique japonaise, « Loi sur le droit d'auteur (partiellement non appliquée) » (Ministère de la Justice, Japon, 2024) ; <https://www.japaneselawtranslation.go.jp/en/laws/view/4207>.
- 863 Ministère israélien de la Justice, « Avis : Utilisations de matériels protégés par le droit d'auteur pour l'apprentissage automatique » (gouvernement israélien, 2022) ; <https://www.gov.il/BlobFolder/legalinfo/machine-learning/he/18-12-2022.pdf>.
- 864 Office de la propriété intellectuelle de Singapour, « Droit d'auteur : fiche d'information sur la loi sur le droit d'auteur de 2021 » (IPOS, 2022) ; <https://www.ipos.gov.sg/docs/default-source/resources-library/copyright/copyright-act-factsheet.pdf>.

- 865 P. Henderson, X. Li, D. Jurafsky, T. Hashimoto, MA Lemley, P. Liang, Modèles de fondation et utilisation équitable, arXiv [cs.CY] (2023) ; <http://arxiv.org/abs/2303.15715>.
- 866 BLW Sobel, La crise de l'utilisation équitable de l'intelligence artificielle. *The Columbia Journal of Law & the Arts* 41, 45–97 (2018) ; <https://doi.org/10.7916/jla.v41i1.2036>.
- 867 MA Lemley, B. Casey, Fair Learning. *Revue de droit du Texas* 99, 743–786 (2020-2021) ; <https://heinonline.org/HOL/P?h=hein.journals/tr99&i=777>.
- 868 P. Samuelson, L'IA générative rencontre le droit d'auteur. *Science* 381, 158–161 (2023) ; <https://doi.org/10.1126/science.adi0656>.
- 869 Tremblay c. OpenAI, Inc. (3:23-cv-03223) Document 1 (2023); https://storage.courtlistener.com/recap/gov.uscourts.cand.414822/gov.uscourts.cand.414822.1.0_1.pdf.
- 870 D. Zhang, B. Xia, Y. Liu, X. Xu, T. Hoang, Z. Xing, M. Staples, Q. Lu, L. Zhu, « Protection de la vie privée et du droit d'auteur dans « IA générative : une perspective du cycle de vie » dans la 3e Conférence internationale sur l'ingénierie de l'IA - Ingénierie logicielle pour l'IA (CAIN) (Lisbonne, Portugal, 2024) ; <http://arxiv.org/abs/2311.18252>.
- 871 R. Mahari, S. Longpre, « Discit Ergo Est : provenance des données de formation et utilisation équitable » dans *Dynamics of Generative AI*, T. Schrepel, V. Stocker, éd. (Network Law Review, 2023) ; <https://www.networklawreview.org/mahari-longpre-generative-ai/>.
- 872 K. Lee, AF Cooper, J. Grimmelmann, « Talkin' 'Bout AI Generation : Copyright and the Generative-IA Supply Chain (The Short Version) » dans *Actes du Symposium sur l'informatique et le droit (CSLAW '24)* (Association for Computing Machinery, New York, NY, États-Unis, 2024), pp. 48–63 ; <https://doi.org/10.1145/3614407.3643696>.
- 873 J. Grimmelmann, Droits d'auteur pour les robots alphabétisés. *Iowa Law Review* 101, 657–682 (2015-2016) ; <https://heinonline.org/HOL/P?h=hein.journals/ilr101&i=681>.
- 874 K. Lee, AF Cooper, J. Grimmelmann, D. Ippolito, IA et droit : la prochaine génération (2023) ; <https://doi.org/10.2139/ssrn.4580739>.
- 875 L. Tiedrich, Lorsque l'IA génère du travail, les clauses contractuelles types peuvent aider à générer de la valeur et de la clarté, *Observatoire des politiques de l'IA de l'OCDE* (2024) ; <https://oecd.ai/en/wonk/contractual-terms>.
- 876 M. Sag, Sécurité du droit d'auteur pour l'IA générative. *Houston Law Review / Université de Houston* 61, 295–347 (2023) ; <https://houstonlawreview.org/article/92126-copyright-safety-for-generative-ai>.
- 877 N. Vyas, SM Kakade, B. Barak, « Sur la protection prouvable du droit d'auteur pour les modèles génératifs » dans *Actes de la 40e Conférence internationale sur l'apprentissage automatique (ICML 2023)* (PMLR, Kigali, Rwanda, 2023) ; <https://proceedings.mlr.press/v202/vyas23b.html>.
- 878 L. Soldaini, R. Kinney, A. Bhagia, D. Schwenk, D. Atkinson, R. Authur, B. Bogin, K. Chandu, J. Dumas, Y. Elazar, V. Hofmann, AH Jha, S. Kumar, L. Lucy, X. Lyu, N. Lambert, I. Magnusson, ... K. Lo, Dolma : Un corpus ouvert de trois mille milliards de jetons pour la recherche sur la préformation des modèles de langage, arXiv [cs.CL] (2024) ; <http://arxiv.org/abs/2402.00159>.
- 879 EM Bender, B. Friedman, Déclarations de données pour le traitement du langage naturel : vers une atténuation des biais du système et une meilleure science. *Transactions of the Association for Computational Linguistics* 6, 587–604 (2018) ; https://doi.org/10.1162/tacl_a_00041.
- 880 R. Bommasani, K. Klyman, S. Longpre, S. Kapoor, N. Maslej, B. Xiong, D. Zhang, P. Liang, « L'indice de transparence du modèle de fondation » (Centre de recherche sur les modèles de fondation (CRFM) et Institut sur l'intelligence artificielle centrée sur l'humain (HAI), 2023) ; <http://arxiv.org/abs/2310.12941>.
- 881 R. Mahari, L. Shayne, L. Donewald, A. Polozov, A. Pentland, A. Lipsitz, Commentaire au US Copyright Office sur les données Provenance et droits d'auteur. *US Copyright Office* (2023) ; <https://dspace.mit.edu/handle/1721.1/154171?show=full?show=full>.
- 882 B. Magagna, D. Goldfarb, P. Martin, M. Atkinson, S. Koulouzis, Z. Zhao, « Provenance des données » dans *Vers des infrastructures de recherche interopérables pour les sciences de l'environnement et de la Terre : une approche guidée par un modèle de référence pour les défis communs*, Z. Zhao, M. Hellström, éd. (Springer International Publishing, Cham, 2020), pp. 208–225 ; https://doi.org/10.1007/978-3-030-52829-4_12.
- 883 S. Longpre, R. Mahari, N. Obeng-Marnu, W. Brannon, T. South, J. Kabbara, S. Pentland, L'authenticité des données, le consentement et la provenance de l'IA sont tous brisés : que faudra-t-il pour les réparer ? Une exploration du MIT sur l'IA générative (2024) ; <https://doi.org/10.21428/e4baedd9.a650f77d>.
- 884 KI Gero, M. Desai, C. Schnitzler, N. Eom, J. Cushman, EL Glassman, Attitudes des écrivains créatifs sur l'écriture en tant que Données de formation pour les grands modèles linguistiques, arXiv [cs.HC] (2024) ; <http://arxiv.org/abs/2409.14281>.
- 885 R. Fletcher, « Combien de sites d'actualités bloquent les robots d'exploration de l'IA ? » (Reuters Institute for the Study of Journalism, 2024) ; <https://doi.org/10.60625/RISJ-XM9G-WS87>.
- 886 Commission européenne, AI Act : Participer à l'élaboration du premier code de bonnes pratiques en matière d'IA à usage général,

- Façonner l'avenir numérique de l'Europe (2024) ; <https://digital-strategy.ec.europa.eu/en/news/ai-act-participate-drawing-first-general-purpose-ai-code-practice> .
- 887 Institut national des normes et de la technologie (NIST), Cadre de gestion des risques de l'IA (2021) ; <https://www.nist.gov/itl/ai-risk-management-framework>.
- 888 J. Lee, T. Le, J. Chen, D. Lee, « Les modèles linguistiques plagient-ils ? » dans Actes de la conférence Web ACM 2023 (ACM, New York, NY, États-Unis, 2023) ; <https://doi.org/10.1145/3543507.3583199>.
- 889 AF Cooper, J. Grimmelmann, Les fichiers sont dans l'ordinateur : sur le droit d'auteur, la mémorisation et l'IA générative. *Revue de droit Chicago-Kent* (2024) ; https://blog.genlaw.org/pdfs/genlaw_icml2024/5.pdf.
- 890 C. Zhang, D. Ippolito, K. Lee, M. Jagielski, F. Tramèr, N. Carlini, « Mémorisation contrefactuelle dans les modèles de langage neuronal » dans 37e Conférence internationale sur les systèmes de traitement de l'information neuronale (NeurIPS 2023) (Curran Associates Inc., Red Hook, NY, États-Unis, 2023) ; <https://dl.acm.org/doi/10.5555/3666122.3667830>.
- 891 L. He, Y. Huang, W. Shi, T. Xie, H. Liu, Y. Wang, L. Zettlemoyer, C. Zhang, D. Chen, P. Henderson, Fantastique Les bêtes protégées par le droit d'auteur et comment (ne pas) les générer, arXiv [cs.CV] (2024) ; <http://arxiv.org/abs/2406.14526>.
- 892 S. Liu, Y. Yao, J. Jia, S. Casper, N. Baracaldo, P. Hase, X. Xu, Y. Yao, H. Li, KR Varshney, M. Bansal, S. Koyejo, Y. Liu, Repenser le désapprentissage automatique pour les grands modèles linguistiques, arXiv [cs.LG] (2024) ; <http://arxiv.org/abs/2402.08787>.
- 893* R. Eldan, M. Russinovich, Qui est Harry Potter ? Désapprentissage approximatif dans les LLM, arXiv [cs.CL] (2023) ; <http://arxiv.org/abs/2310.02238>.
- 894 T. Chen, A. Asai, N. Mireshghallah, S. Min, J. Grimmelmann, Y. Choi, H. Hajishirzi, L. Zettlemoyer, PW Koh, CopyBench : Mesure de la reproduction littérale et non littérale de textes protégés par le droit d'auteur dans la génération de modèles linguistiques, arXiv [cs.CL] (2024) ; <http://arxiv.org/abs/2407.07087>.
- 895 TT Nguyen, TT Huynh, P. Le Nguyen, AW-C. Liew, H. Yin, QVH Nguyen, Une enquête sur le désapprentissage automatique, arXiv [cs.LG] (2022) ; <http://arxiv.org/abs/2209.02299>.
- 896 T. Baumhauer, P. Schöttle, M. Zeppelzauer, Machine Unlearning : Filtration linéaire pour les classificateurs basés sur Logit. *Apprentissage automatique* 111, 3203–3226 (2022) ; <https://doi.org/10.1007/s10994-022-06178-9>.
- 897 Z. Liu, H. Ye, C. Chen, Y. Zheng, K.-Y. Lam, Menaces, attaques et défenses dans le désapprentissage automatique : une enquête, arXiv [cs.CR] (2024) ; <http://arxiv.org/abs/2403.13682>.
- 898 J. Xu, Z. Wu, C. Wang, X. Jia, Machine Unlearning: Solutions and Challenges. *IEEE Transactions on Emerging Topics in Computational Intelligence* 8, 2150–2168 (2024) ; <https://doi.org/10.1109/tetci.2024.3379240>.
- 899 S. Nevo, D. Lahav, A. Karpur, Y. Bar-On, HA Bradley, J. Alstott, Sécurisation des poids des modèles d'IA : prévention du vol et de l'utilisation abusive des modèles Frontier (RAND Corporation, Santa Monica, CA, 2024) ; <https://doi.org/10.7249/RRA2849-1>.
- 900 R. Bommasani, S. Kapoor, K. Klyman, S. Longpre, A. Ramaswami, D. Zhang, M. Schaake, DE Ho, A. Narayanan, P. Liang, Considérations pour la gouvernance des modèles de fondations ouvertes. *Science* (New York, NY) 386, 151–153 (2024) ; <https://doi.org/10.1126/science.adp1848>.
- 901 US National Telecommunications and Information Administration, « Modèles de fondation à double usage avec une large Pondérations des modèles disponibles, rapport NTIA (département du Commerce des États-Unis, 2024) ; <https://www.ntia.gov/issues/artificial-intelligence/open-model-weights-report>.
- 902 E. Seger, N. Dreksler, R. Moulange, E. Dardaman, J. Schuett, K. Wei, C. Winter, M. Arnold, S. Ó. hÉigeartaigh, A. Korinek, M. Anderljung, B. Bucknall, A. Chan, E. Stafford, L. Koessler, A. Ovadya, B. Garfinkel, ... A. Gupta, « Modèles de fondation hautement performants en open source : une évaluation des risques, des avantages et des méthodes alternatives pour poursuivre les objectifs de l'open source » (Centre pour la gouvernance de l'IA, 2023) ; <http://arxiv.org/abs/2311.09227>.
- 903 P. Gade, S. Lermen, C. Rogers-Smith, J. Ladish, BadLlama : Supprimer à moindre coût les réglages de sécurité de Llama 2-Discussion 13B, arXiv [cs.CL] (2023) ; <http://arxiv.org/abs/2311.00117>.
- 904* A. Zou, Z. Wang, N. Carlini, M. Nasr, J. Zico Kolter, M. Fredrikson, Attaques adverses universelles et transférables contre Modèles de langage alignés, arXiv [cs.CL] (2023) ; <http://arxiv.org/abs/2307.15043>.
- 905 I. Yum, Agents linguistiques et conception malveillante. *Philosophie et technologie* 37, 1–19 (2024) ; <https://doi.org/10.1007/s13347-024-00794-0>.
- 906 S. Lermen, C. Rogers-Smith, J. Ladish, Le réglage fin de LoRA annule efficacement la formation à la sécurité dans Llama 2-Chat 70B, arXiv [cs.LG] (2023) ; <http://arxiv.org/abs/2310.20624>.
- 907 A. Arditi, O. Obeso, A. Syed, D. Paleka, N. Panickssery, W. Gurnee, N. Nanda, Le refus dans les modèles linguistiques est médiatisé par une seule direction, arXiv [cs.LG] (2024) ; <http://arxiv.org/abs/2406.11717>.
- 908 J. Cable, A. Black, « Avec l'intelligence artificielle open source, n'oubliez pas les leçons des logiciels open source » (Agence de cybersécurité et de sécurité des infrastructures CISA, 2024) ; <https://www.cisa.gov/news->

événements/actualités/intelligence-artificielle-open-source-n'oubliez-pas-les-leçons-logiciels-open-source.

- 909 par J. Bateman, D. Baer, SA Bell, GO Brown, M.-F. (tino) Cuéllar, D. Ganguli, P. Henderson, B. Kotila, L. Lessig, N.-B. François Lundblad, J. Napolitano, D. Raji, E. Seger, M. Sheehan, A. Skowron, I. Solaiman, H. Toner, AP Zvyagina, « Au-delà de l'ouverture et de la fermeture : consensus émergent et questions clés pour la gouvernance du modèle d'IA de base » (Carnegie Endowment for International Peace, 2024) ; <https://carnegieendowment.org/research/2024/07/beyond-open-vs-closed-emerging-consensus-and-key-questions-for-foundation-ai-model-governance?lang=en> .
- 910 E. Seger, B. O'Dell, « Open Horizons : Exploration des approches techniques et politiques nuancées de l'ouverture dans l'IA » (Demos, 2024) ; <https://demos.co.uk/research/open-horizons-exploring-nuanced-technical-and-policy-approachs-to-openness-in-ai/> .
- 911 François S. Kapoor, R. Bommasani, K. Klyman, S. Longpre, A. Ramaswami, P. Cihon, AK Hopkins, K. Bankston, S. Biderman, M. Bogen, R. Chowdhury, A. Engler, P. Henderson, Y. Jernite, S. Lazar, S. Maffulli, A. Nelson, ... A. Narayanan, « Position : sur l'impact sociétal des modèles de fondations ouvertes » dans Conférence internationale sur l'apprentissage automatique (PMLR, 2024), pp. 23082–23104 ; <https://proceedings.mlr.press/v235/kapoor24a.html> .
- 912* S. Lakatos, « Une image révélatrice : les images de « déshabillage » générées par l'IA passent des forums de discussion sur la pornographie de niche à une entreprise en ligne à grande échelle et monétisée » (Graphika, 2023) ; <https://graphika.com/reports/a-reveling-picture> .
- 913 D. Thiel, M. Stroebel, R. Portnoff, « ML génératif et CSAM : implications et atténuations » (Thorn & Stanford Observatoire Internet, 2023) ; <https://fsi.stanford.edu/publication/generative-ml-and-csam-implications-and-mitigations> .
- 914 A. Engler, « Comment les logiciels open source façonnent la politique de l'IA » (Brookings, 2021) ; <https://www.brookings.edu/articles/how-open-source-software-shapes-ai-policy/> .
- 915 D. Gray Widder, S. West, M. Whittaker, Open (for Business) : Big Tech, pouvoir concentré et économie politique de l'IA ouverte, SSRN [préimpression] (2023) ; <https://doi.org/10.2139/ssrn.4543807> .
- 916 K. Blind, M. Böhm, P. Grzegorzewska, A. Katz, S. Muto, S. Pätsch, T. Schubert, « Étude sur l'impact des logiciels et du matériel open source sur l'indépendance technologique, la compétitivité et l'innovation dans l'économie de l'UE, rapport d'étude final » (Commission européenne, 2021) ; <https://digital-strategy.ec.europa.eu/en/library/study-about-impact-open-source-software-and-hardware-technological-independence-competitiveness-and> .
- 917 Y. Kilcher, Ykilcher/gpt-4chan (2023) ; <https://huggingface.co/ykilcher/gpt-4chan> .
- 918 S. Borgeaud, A. Mensch, J. Hoffmann, T. Cai, E. Rutherford, K. Millican, G. van den Driessche, J.-B. Lespiau, B. Damoc, A. Clark, D. de Las Casas, A. Guy, J. Menick, R. Ring, T. Hennigan, S. Huang, L. Maggiore, ... L. Sifre, Améliorer les modèles de langage en récupérant des milliards de jetons. Conférence internationale sur l'apprentissage automatique 162, 2206–2240 (2021) ; <https://proceedings.mlr.press/v162/borgeaud22a/borgeaud22a.pdf> .
- 919 P. Henderson, E. Mitchell, C. Manning, D. Jurafsky, C. Finn, « Modèles autodestructeurs : augmentation des coûts de « Utilisations doubles néfastes des modèles de fondation » dans les actes de la conférence 2023 de l'AAAI/ACM sur l'IA, l'éthique et la société (Association for Computing Machinery, New York, NY, États-Unis, 2023) AIES '23, pp. 287–296 ; <https://doi.org/10.1145/3600211.3604690> .
- 920 J. Deng, S. Pang, Y. Chen, L. Xia, Y. Bai, H. Weng, W. Xu, SOPHON : Apprentissage non affiné pour restreindre la tâche Transférabilité pour les modèles pré-entraînés, arXiv [cs.LG] (2024) ; <http://arxiv.org/abs/2404.12699> .
- 921 T. Huang, S. Hu, L. Liu, « Vaccin : Alignement sensible aux perturbations pour les grands modèles de langage contre les attaques de réglage fin nuisibles » dans 38e conférence annuelle sur les systèmes de traitement de l'information neuronale (NeurIPS 2024) (2024) ; <https://openreview.net/pdf?id=lpXDZKiAnt> .
- 922 D. Rosati, J. Wehner, K. Williams, Ł. Bartoszcze, D. Atanasov, R. Gonzales, S. Majumdar, C. Maple, H. Sajjad, F. Rudzicz, Le bruit de représentation empêche efficacement les réglages fins nuisibles sur les LLM, arXiv [cs.CL] (2024) ; <http://arxiv.org/abs/2405.14577> .
- 923 R. Tamirisa, B. Bharathi, L. Phan, A. Zhou, A. Gatti, T. Suresh, M. Lin, J. Wang, R. Wang, R. Arel, A. Zou, D. Song, B. Li, D. Hendrycks, M. Mazeika, Garanties inviolables pour les LLM à pondération ouverte, arXiv [cs.LG] (2024) ; <http://arxiv.org/abs/2408.00761> .
- 924 G. Wang, Y.-N. Chuang, R. Tang, S. Zhong, J. Yuan, H. Jin, Z. Liu, V. Chaudhary, S. Xu, J. Caverlee, X. Hu, Taylor Unswift : libération de poids sécurisée pour les grands modèles de langage via l'expansion de Taylor, arXiv [cs.CR] (2024) ; <http://arxiv.org/abs/2410.05331> .
- 925 M. Srikumar, J. Chang, K. Chmielinski, « Stratégies d'atténuation des risques pour la chaîne de valeur du modèle de fondation ouverte : « Résultats de l'atelier PAI co-organisé avec GitHub » (Partenariat sur l'IA, 2024) ; https://partnershiponai.org/wp-content/uploads/dlm_uploads/2024/07/open-foundation-model-risk-mitigation_rev3-1.pdf .

- 926 E. David, Meta dévoile son modèle d'IA le plus puissant, Llama 3.1, avec 405 milliards de paramètres, VentureBeat (2024) ; <https://venturebeat.com/ai/meta-unleashes-its-most-powerful-ai-model-llama-3-1-with-405b-parameters/>.
- 927 B. Muralidharan, H. Beadles, R. Marzban, KS Mupparaju, Knowledge AI : Affiner les modèles NLP pour faciliter l'extraction et la compréhension des connaissances scientifiques, arXiv [cs.CL] (2024) ; <http://arxiv.org/abs/2408.04651>.
- 928* L. Weidinger, M. Rauh, N. Marchal, A. Manzini, LA Hendricks, J. Mateos-Garcia, S. Bergman, J. Kay, C. Griffin, B. Bariach, I. Gabriel, V. Rieser, W. Isaac, « Évaluation sociotechnique de la sécurité des systèmes d'IA générative » (Google Deepmind, 2023) ; <http://arxiv.org/abs/2310.11986>.
- 929* L. Weidinger, J. Barnhart, J. Brennan, C. Butterfield, S. Young, W. Hawkins, L. A. Hendricks, R. Comanescu, O. Chang, M. Rodriguez, J. Beroshi, D. Bloxwich, L. Proleev, J. Chen, S. Farquhar, L. Ho, I. Gabriel, ... W. Isaac, « Évaluations holistiques de la sécurité et de la responsabilité des modèles d'IA avancés » (Google Deepmind, 2024) ; <http://arxiv.org/abs/2404.14068>.
- 930* I. Solaiman, Z. Talat, W. Agnew, L. Ahmad, D. Baker, S. L. Blodgett, H. Daumé III, J. Dodge, E. Evans, S. Hooker, Y. Jernite, AS Luccioni, A. Lusoli, M. Mitchell, J. Newman, M.-T. Png, A. Strait, A. Vassilev, Évaluation de l'impact social des systèmes d'IA générative dans les systèmes et la société, arXiv [cs.CY] (2023) ; <http://arxiv.org/abs/2306.05949>.
- 931 ARR Salammagari, G. Srivastava, Faire progresser la compréhension du langage naturel pour les langues à faibles ressources : Progrès actuels, applications et défis. *Revue internationale de recherche avancée en ingénierie et technologie* 15, 244–255 (2024) ; https://iaeme.com/Home/article_id/IJARET_15_03_021.
- 932 A. Birhane, W. Isaac, V. Prabhakaran, M. Diaz, MC Elish, I. Gabriel, S. Mohamed, « Le pouvoir au peuple ? « Opportunités et défis pour l'IA participative » dans les actes de la 2e conférence de l'ACM sur l'équité et l'accès aux algorithmes, mécanismes et optimisation (EAAMO '22) (Association for Computing Machinery, New York, NY, États-Unis, 2022), pp. 1–8 ; <https://doi.org/10.1145/3551624.3555290>.
- 933 P. Slattery, AK Saeri, EAC Grundy, J. Graham, M. Noetel, R. Uuk, J. Dao, S. Pour, S. Casper, N. Thompson, The AI Risk Repository : une méta-analyse complète, une base de données et une taxonomie des risques liés à l'intelligence artificielle, arXiv [cs.AI] (2024) ; <http://arxiv.org/abs/2408.12622>.
- 934 Partenariat sur l'IA, « [Projet] Lignes directrices pour une IA participative et inclusive » (2024) ; <https://partnershiponai.notion.site/1e8a6131dda045f1ad00054933b0bda0?v=dc890146f7d464a86f11fcd5de372c0>.
- 935 M. Maghsoudi, A. Mohammadi, S. Habibipour, Naviguer et répondre aux préoccupations du public en matière d'IA : enseignements tirés de l'analyse des médias sociaux et de Delphi. *IEEE Access : Practical Innovations, Open Solutions* 12, 1–1 (2024) ; <https://doi.org/10.1109/access.2024.3440660>.
- 936 K. Grosse, L. Bieringer, TR Besold, AM Alahi, « Vers des modèles de menaces plus pratiques dans la sécurité de l'intelligence artificielle » dans 33e Symposium sur la sécurité USENIX (USENIX Security 24) (2024), pp. 4891–4908 ; <https://www.usenix.org/system/files/usenixsecurity24-grosse.pdf>.
- 937 H. Li, Z. Ren, M. Fan, W. Li, Y. Xu, Y. Jiang, W. Xia, Examen des méthodes d'analyse de scénarios dans la planification et l'exploitation des systèmes électriques modernes : méthodologies, applications et défis. *Electric Power Systems Research* 205, 107722 (2022) ; <https://doi.org/10.1016/j.epr.2021.107722>.
- 938 A. Mantelero, L'évaluation de l'impact sur les droits fondamentaux (FRIA) dans la loi sur l'IA : racines, obligations légales et éléments clés pour un modèle. *Computer Law and Security Report* 54, 106020 (2024) ; <https://doi.org/10.1016/j.clsr.2024.106020>.
- 939 ID Raji, P. Xu, C. Honigsberg, D. Ho, « Surveillance externe : conception d'un écosystème d'audit tiers pour la gouvernance de l'IA » dans Actes de la conférence 2022 de l'AAAI/ACM sur l'IA, l'éthique et la société (AIES '22) (Association for Computing Machinery, New York, NY, États-Unis, 2022), pp. 557–571 ; <https://doi.org/10.1145/3514094.3534181>.
- 940 V. Storchan, R. Kumar, R. Chowdhury, S. Goldfarb-Tarrant, S. Cattell, « 2024 Générative AI Red Teaming Rapport de transparence (Humane intelligence, 2024).
- 941* S. Wan, C. Nikolaidis, D. Song, D. Molnar, J. Crnkovich, J. Grace, M. Bhatt, S. Chennabasappa, S. Whitman, S. Ding, V. Ionescu, Y. Li, J. Saxe, CYBERSECEVAL 3 : Faire progresser l'évaluation des risques et des capacités de cybersécurité dans les grands modèles linguistiques, arXiv [cs.CR] (2024) ; <http://arxiv.org/abs/2408.01605>.
- 942 RJ Neuwirth, Pratiques d'intelligence artificielle interdites dans la proposition de loi sur l'intelligence artificielle de l'UE (AIA). *Revue de droit et de sécurité informatique* 48, 105798 (2023) ; <https://doi.org/10.1016/j.clsr.2023.105798>.
- 943 L. Heim, L. Koessler, Seuils de calcul d'entraînement : caractéristiques et fonctions dans la régulation de l'IA, arXiv [cs.CY] (2024) ; <http://arxiv.org/abs/2405.10799>.
- 944 L. Koessler, J. Schuett, M. Anderljung, Seuils de risque pour l'IA de pointe, arXiv [cs.CY] (2024) ; <http://arxiv.org/abs/2406.14713>.

- 945 Centre pour la sécurité des procédés chimiques (CCPS), Nœuds papillon dans la gestion des risques (John Wiley & Sons, Nashville, TN, 2018) ; <https://doi.org/10.1002/9781119490357>.
- 946 Organisation internationale de normalisation, « ISO 21448:2022 : Véhicules routiers — Sécurité du véhicule prévu « Fonctionnalité » (ISO, 2022) ; <https://www.iso.org/standard/77490.html>.
- 947* Politique de mise à l'échelle anthropique et responsable. (2024) ; <https://assets.anthropic.com/m/24a47b00f10301cd/original/Anthropic-Responsible-Scaling-Policy-2024-10-15.pdf>.
- Partenariat 948 sur l'IA, Guide du PAI pour un déploiement sûr du modèle de fondation (2023) ; <https://partnershiponai.org/modeldeployment/>.
- 949 T. Kelly, Une approche systématique de la gestion des dossiers de sécurité. SAE Transactions : Journal of Materials & Manufacturing 113, 257–266 (2004) ; <http://www.jstor.org/stable/44699541>.
- 950 B. Lakshmi Prasanna, M. SaidiReddy, (CSM2-RA-R2-TI) : Modèle de maturité de la cybersécurité pour l'évaluation des risques à l'aide d'un registre des risques pour le renseignement sur les menaces. Journal of Physics. Conference Series 2040, 012005 (2021) ; <https://doi.org/10.1088/1742-6596/2040/1/012005>.
- 951* Y. Zeng, K. Klyman, A. Zhou, Y. Yang, M. Pan, R. Jia, D. Song, P. Liang, B. Li, AI Risk Categorization Decoded (AIR 2024) : des réglementations gouvernementales aux politiques d'entreprise, arXiv [cs.CY] (2024) ; <http://arxiv.org/abs/2406.17864>.
- 952 H. Wu, AI Whistleblowers, SSRN [préimpression] (2024) ; <https://doi.org/10.2139/ssrn.4790511>.
- 953 MITRE ATLAS, Incidents d'IA MITRE ATLAS (2024) ; <https://ai-incidents.mitre.org/>.
- 954 B. Robinson, J. Ginns, « Transformer la gouvernance des risques dans les entreprises de pointe en IA » (Centre pour la résilience à long terme, 2024) ; <https://www.longtermresilience.org/wp-content/uploads/2024/07/Transforming-risk-government-at-frontier-AI-companies-CLTR-1.pdf>.
- 955 J. Schuett, Trois lignes de défense contre les risques liés à l'IA. IA et société (2023) ; <https://doi.org/10.1007/s00146-023-01811-0>.
- 956 R. Bommasani, K. Klyman, S. Longpre, B. Xiong, S. Kapoor, N. Maslej, A. Narayanan, P. Liang, Modèle de fondation Rapports de transparence, arXiv [cs.LG] (2024) ; <http://arxiv.org/abs/2402.16268>.
- 957* D. Hendrycks, N. Carlini, J. Schulman, J. Steinhardt, Problèmes non résolus en matière de sécurité du ML, arXiv [cs.LG] (2021) ; <http://arxiv.org/abs/2109.13916>.
- 958 M. Anderjung, ET Smith, J. O'Brien, L. Soder, B. Bucknall, E. Bluemke, J. Schuett, R. Trager, L. Strahm, R. Chowdhury, Vers des LLM de pointe publiquement responsables : construire un écosystème de contrôle externe dans le cadre d'ASPIRE, arXiv [cs.CY] (2023) ; <http://arxiv.org/abs/2311.14711>.
- 959 R. Gupta, L. Walker, R. Corona, S. Fu, S. Petryk, J. Napolitano, T. Darrell, AW Reddie, Gouvernance de l'IA centrée sur les données : aborder les limites des politiques axées sur les modèles, arXiv [cs.CY] (2024) ; <http://arxiv.org/abs/2409.17216>.
- 960 D. McDuff, T. Korjakow, S. Cambo, JJ Benjamin, J. Lee, Y. Jernite, CM Ferrandis, A. Gokaslan, A. Tarkowski, J. Lindley, AF Cooper, D. Contractor, Sur la normalisation des clauses d'utilisation comportementale et leur adoption pour une licence responsable de l'IA, arXiv [cs.SE] (2024) ; <http://arxiv.org/abs/2402.05979>.
- 961 B. Rakova, J. Yang, H. Cramer, R. Chowdhury, Où l'IA responsable rencontre la réalité : perspectives des praticiens sur Facteurs favorisant l'évolution des pratiques organisationnelles. Actes de l'ACM sur l'interaction homme-machine 5, 1–23 (2021) ; <https://doi.org/10.1145/3449081>.
- 962* Microsoft AI, « Mettre les principes en pratique : comment nous abordons l'IA responsable chez Microsoft » (Microsoft, 2020) ; <https://www.microsoft.com/cms/api/am/binary/RE4pKH5>.
- 963 J. Schuett, A.-K. Reuel, A. Carlier, Comment concevoir un comité d'éthique de l'IA. AI and Ethics, 1–19 (2024) ; <https://doi.org/10.1007/s43681-023-00409-y>.
- 964 G. de Beco, Évaluations d'impact sur les droits de l'homme. Revue trimestrielle néerlandaise des droits de l'homme 27, 139–166 (2009) ; <https://doi.org/10.1177/016934410902700202>.
- 965 E. Donahoe, MM Metzger, Intelligence artificielle et droits de l'homme. Journal of Democracy 30, 115–126 (2019) ; <https://doi.org/10.1353/jod.2019.0029>.
- 966 S. Makridakis, L'art et la science de la prévision Une évaluation et des orientations futures. Revue internationale de Prévisions 2, 15–39 (1986) ; [https://doi.org/10.1016/0169-2070\(86\)90028-2](https://doi.org/10.1016/0169-2070(86)90028-2).
- 967 E. Karger, P. Atanasov, PE Tetlock, « Améliorer les jugements sur le risque existentiel : meilleures prévisions, questions, explications, politiques » (Future of Humanity Institute, 2022) ; <https://www.fhi.ox.ac.uk/wp-content/uploads/2022/05/Improving-Judgments-of-Existential-Risk.pdf>.
- 968 L. Koessler, J. Schuett, Évaluation des risques dans les entreprises AGI : un aperçu des techniques courantes d'évaluation des risques

- d'autres industries critiques pour la sécurité, arXiv [cs.CY] (2023) ; <http://arxiv.org/abs/2307.08823>.
- 969 B. Anderson-Samways, « Précédents réglementaires pertinents pour l'IA : une recherche systématique dans toutes les agences fédérales » (Institut pour la politique et la stratégie de l'IA, 2024) ; <https://www.iaps.ai/research/ai-relevant-regulatory-precedent>.
- 970 HE Roland, B. Moriarty, Ingénierie et gestion de la sécurité des systèmes (Wiley, New York, 2e éd., 1990) ; <https://www.wiley.com/en-us/System+Safety+Engineering+and+Management%2C+2nd+Edition-p-9780471618164>.
- 971 NG Leveson, Ingénierie d'un monde plus sûr : la pensée systémique appliquée à la sécurité (The MIT Press, 2012) ; <https://doi.org/10.7551/mitpress/8179.001.0001>.
- 972 S. Dekker, Fondements de la science de la sécurité : Un siècle de compréhension des accidents et des catastrophes (Routledge, Londres, Angleterre, 2019) ; <https://doi.org/10.4324/9781351059794>.
- 973 ISO, ISO 31000 : Gestion des risques, ISO (2018) ; <https://www.iso.org/iso-31000-gestion-des-risques.html>.
- 974 E. Black, R. Naidu, R. Ghani, K. Rodolfa, D. Ho, H. Heidari, « Vers une opérationnalisation de l'équité du ML tenant compte des pipelines : un programme de recherche pour le développement de lignes directrices et d'outils pratiques » dans Actes de la 3e conférence de l'ACM sur l'équité et l'accès aux algorithmes, mécanismes et optimisation (EAAMO '23) (Association for Computing Machinery, New York, NY, États-Unis, 2023), pp. 1–11 ; <https://doi.org/10.1145/3617694.3623259>.
- 975 S. Rismani, R. Shelby, A. Smart, E. Jatho, J. Kroll, A. Moon, N. Rostamzadeh, « Des accidents d'avion aux dommages algorithmiques : applicabilité des cadres d'ingénierie de sécurité pour un ML responsable » dans Actes de la conférence CHI 2023 sur les facteurs humains dans les systèmes informatiques (CHI '23) (Association for Computing Machinery, New York, NY, États-Unis, 2023), pp. 1–18 ; <https://doi.org/10.1145/3544548.3581407>.
- 976 R. Hawkins, C. Paterson, C. Picardi, Y. Jia, R. Calinescu, I. Habli, Orientations sur l'assurance de l'apprentissage automatique dans Systèmes autonomes (AMLAS), arXiv [cs.LG] (2021) ; <http://arxiv.org/abs/2102.01564>.
- 977 T. Raz, D. Hillson, Une étude comparative des normes de gestion des risques. Gestion des risques : une norme internationale Journal 7, 53–66 (2005) ; <https://doi.org/10.1057/palgrave.rm.8240227>.
- 978 J. Clymer, N. Gabrieli, D. Krueger, T. Larsen, Cas de sécurité : comment justifier la sécurité des systèmes d'IA avancés, arXiv [cs.CY] (2024) ; <http://arxiv.org/abs/2403.10462>.
- 979 C. Haddon-Cave, The Nimrod Review : Un examen indépendant des questions plus vastes entourant la perte de l'avion RAF Nimrod MR2 XV230 en Afghanistan en 2006, rapport (Stationery Office, 2009) ; <https://www.gov.uk/government/publications/the-nimrod-review>.
- 980 NG Leveson, Application de la pensée systémique à l'analyse et à l'apprentissage des événements. Safety Science 49, 55–64 (2011) ; <https://doi.org/10.1016/j.ssci.2009.12.021>.
- 981 D. Hendrycks, Introduction à la sécurité, à l'éthique et à la société de l'IA (CRC Press, 2024) ; <https://www.aisafetybook.com/>.
- 982 O. Delaney, O. Guest, Z. Williams, Cartographie de la recherche sur la sécurité technique dans les entreprises d'IA : une revue de la littérature et Analyse des incitations, arXiv [cs.CY] (2024) ; <http://arxiv.org/abs/2409.07878>.
- 983 R. Uuk, A. Brouwer, N. Dreksler, V. Pulignano, R. Bommasani, Mesures efficaces d'atténuation des risques systémiques IA à usage général. (2024) ; https://papers.ssrn.com/sol3/papers.cfm?abstract_id=5021463.
- 984 DA Boiko, R. MacKnight, G. Gomes, Capacités de recherche scientifique autonomes émergentes du langage large Modèles, arXiv [physics.chem-ph] (2023) ; <http://arxiv.org/abs/2304.05332>.
- 985 Q. Lu, L. Zhu, X. Xu, Z. Xing, S. Harrer, J. Whittle, Vers une IA générative responsable : une architecture de référence pour la conception d'agents basés sur des modèles de base, arXiv [cs.AI] (2023) ; <http://arxiv.org/abs/2311.13148>.
- 986* Équipe SIMA, M. A. Raad, A. Ahuja, C. Barros, F. Besse, A. Bolt, A. Bolton, B. Brownfield, G. Buttimore, M. Cant, S. Chakera, SCY Chan, J. Clune, A. Collister, V. Copeman, A. Cullum, I. Dasgupta, ... N. Young, « Mise à l'échelle d'agents instructibles dans de nombreux mondes simulés » (Google Deepmind, 2024) ; <http://arxiv.org/abs/2404.10179>.
- 987 T. Schick, J. Dwivedi-Yu, R. Dessi, R. Raileanu, M. Lomeli, E. Hambro, L. Zettlemoyer, N. Cancedda, T. Scialom, « Toolformer : les modèles de langage peuvent apprendre à utiliser des outils » dans la 37e conférence sur les systèmes de traitement de l'information neuronale (NeurIPS 2023) (La Nouvelle-Orléans, LA, États-Unis, 2023) ; <https://openreview.net/forum?id=Yacmpz84TH>.
- 988 Y. Tian, X. Yang, J. Zhang, Y. Dong, H. Su, Evil Geniuses : Exploration de la sécurité des agents basés sur LLM, arXiv [cs.CL] (2023) ; <http://arxiv.org/abs/2311.11855>.
- 989 Z. Wu, C. Han, Z. Ding, Z. Weng, Z. Liu, S. Yao, T. Yu, L. Kong, OS-Copilot : vers des agents informatiques généralistes avec Auto-amélioration, arXiv [cs.AI] (2024) ; <http://arxiv.org/abs/2402.07456>.
- 990 Z. Xi, W. Chen, X. Guo, W. He, Y. Ding, B. Hong, M. Zhang, J. Wang, S. Jin, E. Zhou, R. Zheng, X. Fan, X. Wang, L. Xiong, Y. Zhou, W. Wang, C. Jiang, ... T. Gui, L'essor et le potentiel des agents basés sur des modèles de langage de grande taille : une enquête, arXiv [cs.AI] (2023) ; <http://arxiv.org/abs/2309.07864>.

- 991* T. Masterman, S. Besen, M. Sawtell, A. Chao, Le paysage des architectures émergentes d'agents d'IA pour le raisonnement, la planification et l'appel d'outils : une enquête, arXiv [cs.AI] (2024) ; <http://arxiv.org/abs/2404.11584>.
- 992 M. Hartmann, A. Koller, Une enquête sur les tâches complexes pour les agents interactifs dirigés par des objectifs, arXiv [cs.CL] (2024) ; <http://arxiv.org/abs/2409.18538>.
- 993 T. Xie, D. Zhang, J. Chen, X. Li, S. Zhao, R. Cao, T.J. Hua, Z. Cheng, D. Shin, F. Lei, Y. Liu, Y. Xu, S. Zhou, S. Savarese, C. Xiong, V. Zhong, T. Yu, OSWorld : Analyse comparative des agents multimodaux pour les tâches ouvertes dans des environnements informatiques réels, arXiv [cs.AI] (2024) ; <http://arxiv.org/abs/2404.07972>.
- 994* A. Fourney, G. Bansal, H. Mozannar, C. Tan, E. Salinas, E. (eric) Zhu, F. Niedtner, G. Proebsting, G. Bassman, J. Gerrits, J. Alber, P. Chang, R. Loynd, R. West, V. Dibia, A. Awadallah, E. Kamar, ... S. Amershi, « Magentic-One : un système multi-agent généraliste pour résoudre des tâches complexes » (Microsoft, 2024) ; <https://www.microsoft.com/en-us/research/publication/magentic-one-a-generalist-multi-agent-system-for-solving-complex-tasks/>.
- 995 S. Hu, M. Ouyang, D. Gao, MZ Shou, L'aube de l'agent GUI : une étude de cas préliminaire avec l'ordinateur Claude 3.5 Utilisation, arXiv [cs.AI] (2024) ; <http://arxiv.org/abs/2411.10323>.
- 996 J.-P. Rivera, G. Mukobi, A. Reuel, M. Lamparth, C. Smith, J. Schneider, « Risques d'escalade liés aux modèles linguistiques dans la prise de décision militaire et diplomatique » dans la conférence 2024 de l'ACM sur l'équité, la responsabilité et la transparence (ACM, New York, NY, États-Unis, 2024) ; <https://doi.org/10.1145/3630106.3658942>.
- 997 B. Zhang, Y. Tan, Y. Shen, A. Salem, M. Backes, S. Zannettou, Y. Zhang, Breaking Agents : Compromettre les agents LLM autonomes par l'amplification des dysfonctionnements, arXiv [cs.CR] (2024) ; <http://arxiv.org/abs/2407.20859>.
- 998 K. Greshake, S. Abdelnabi, S. Mishra, C. Endres, T. Holz, M. Fritz, « Pas ce pour quoi vous avez signé : compromettre les applications intégrées au LLM du monde réel avec l'injection indirecte rapide » dans Actes du 16e atelier ACM sur l'intelligence artificielle et la sécurité (AISec '23) (Association for Computing Machinery, New York, NY, États-Unis, 2023), pp. 79–90 ; <https://doi.org/10.1145/3605764.3623985>.
- 999 R. Fang, D. Bowman, D. Kang, Les agents d'IA à commande vocale peuvent réaliser des escroqueries courantes, arXiv [cs.AI] (2024) ; <http://arxiv.org/abs/2410.15650>.
- 1000 M. Andriushchenko, A. Souly, M. Dziemian, D. Duenas, M. Lin, J. Wang, D. Hendrycks, A. Zou, Z. Kolter, M. Fredrikson, E. Winsor, J. Wynne, Y. Gal, X. Davies, AgentHarm : une référence pour mesurer la nocivité des agents LLM, arXiv [cs.LG] (2024) ; <http://arxiv.org/abs/2410.09024>.
- 1001* P. Kumar, E. Lau, S. Vijayakumar, T. Trinh, Scale Red Team, E. Chang, V. Robinson, S. Hendryx, S. Zhou, M. Fredrikson, S. Yue, Z. Wang, Les LLM formés au refus sont facilement jailbreakés en tant qu'agents de navigateur, arXiv [cs.CR] (2024) ; <http://arxiv.org/abs/2410.13886>.
- 1002 A. Chan, C. Ezell, M. Kaufmann, K. Wei, L. Hammond, H. Bradley, E. Bluemke, N. Rajkumar, D. Krueger, N. Kolt, L. Heim, M. Anderljung, Visibilité dans les agents d'IA, arXiv [cs.CY] (2024) ; <http://arxiv.org/abs/2401.13138>.
- 1003 MK Cohen, N. Kolt, Y. Bengio, GK Hadfield, S. Russell, Régulation des agents artificiels avancés. *Science* 384, 36–38 (2024) ; <https://doi.org/10.1126/science.adl0625>.
- 1004 G. Mialon, C. Fourrier, T. Wolf, Y. LeCun, T. Scialom, « GAIA : une référence pour les assistants d'IA généraux » dans la 12e Conférence internationale sur les représentations d'apprentissage (ICLR 2024) (Vienne, Autriche, 2024) ; <https://openreview.net/forum?id=fibxvahvs3>.
- 1005 K. Valmeekam, K. Stechly, S. Kambhampati, « Les LLM ne peuvent toujours pas planifier ; les LRM le peuvent-ils ? Une évaluation préliminaire de l'OpenAI o1 sur PlanBench » dans l'atelier NeurIPS 2024 sur les agents du monde ouvert (2024) ; <https://openreview.net/forum?id=Gcr1Lx4Koz>.
- 1006 PP Liang, A. Zadeh, L.-P. Morency, Fondements et tendances de l'apprentissage automatique multimodal : principes, Défis et questions ouvertes. *ACM Computing Surveys* 56, 1–42 (2024) ; <https://doi.org/10.1145/3656580>.
- 1007 R. Wang, X. Ma, H. Zhou, C. Ji, G. Ye, Y.-G. Jiang, « Jailbreaks multimodaux en boîte blanche contre de grands modèles vision-langage » dans ACM Multimedia 2024 (2024) ; <https://openreview.net/forum?id=SMOQtEaAf>.
- 1008 M. Thiemann, J. Lepoutre, Stitched on the Edge : Évasion des règles, régulateurs intégrés et évolution de Marchés. *American Journal of Sociology* 122, 1775–1821 (2017) ; <https://doi.org/10.1086/691348>.
- 1009 R. Huben, H. Cunningham, LR Smith, A. Ewart, L. Sharkey, « Les autoencodeurs clairsemés trouvent des données hautement interprétables « Caractéristiques des modèles linguistiques » dans la 12e Conférence internationale sur les représentations d'apprentissage (ICLR 2024) (Vienne, Autriche, 2023) ; <https://openreview.net/forum?id=F76bwRSLeK>.
- 1010* L. Gao, TD la Tour, H. Tillman, G. Goh, R. Troll, A. Radford, I. Sutskever, J. Leike, J. Wu, Mise à l'échelle et évaluation des autoencodeurs clairsemés, arXiv [cs.LG] (2024) ; <http://arxiv.org/abs/2406.04093>.
- 1011* T. Lieberum, S. Rajamanoharan, A. Conmy, L. Smith, N. Sonnerat, V. Varma, J. Kramar, A. Dragan, R. Shah, N. Nanda, « Gemma Scope : ouvrez les autoencodeurs clairsemés partout en même temps sur Gemma 2 » dans The 7th BlackboxNLP

- Atelier (2024) ; <https://openreview.net/forum?id=XkMrWOJhNd>.
- 1012 A. Templeton, T. Conerly, J. Marcus, J. Lindsey, T. Bricken, B. Chen, A. Pearce, C. Citro, E. Ameisen, A. Jones, H. Cunningham, NL Turner, C. McDougall, M. MacDiarmid, CD Freeman, TR Summers, E. Rees, ... T. Henighan, Mise à l'échelle de la monosémantique : extraction de caractéristiques interprétables du sonnet de Claude 3. Fil de discussion sur les circuits de transformateurs (2024) ; <https://transformer-circuits.pub/2024/scaling-monosemanticity/index.html>.
- 1013* T. Bricken, A. Templeton, J. Batson, B. Chen, A. Jermyn, T. Conerly, N. Turner, C. Anil, C. Denison, A. Askell, R. Lasenby, Y. Wu, S. Kravec, N. Schiefer, T. Maxwell, N. Joseph, Z. Hatfield-Dodds, ... C. Olah, Vers la monosémantique : décomposition des modèles de langage avec l'apprentissage par dictionnaire, Transformer Circuits Thread (2023) ; <https://transformer-circuits.pub/2023/monosemantic-features>.
- 1014 M. Ananny, K. Crawford, Voir sans savoir : les limites de l'idéal de transparence et son application à Responsabilité algorithmique. *Nouveaux médias et société* 20, 973–989 (2018) ; <https://doi.org/10.1177/1461444816676645>.
- 1015* T. Bolukbasi, A. Pearce, A. Yuan, A. Coenen, E. Reif, F. Viégas, M. Wattenberg, Une illusion d'interprétabilité pour BERT, arXiv [cs.CL] (2021) ; <http://arxiv.org/abs/2104.07143>.
- 1016 K. Kaye, P. Dixon, « Analyse des risques : évaluer et améliorer les outils de gouvernance de l'IA. Une revue internationale des outils de gouvernance de l'IA et des suggestions de voies à suivre » (World Privacy Forum, 2023) ; https://www.worldprivacyforum.org/wp-content/uploads/2023/12/WPF_Risky_Analysis_December_2023_fs.pdf.
- 1017 A. Makelov, G. Lange, A. Geiger, N. Nanda, « Est-ce le sous-espace que vous recherchez ? Une illusion d'interprétabilité pour le patch d'activation du sous-espace » dans la 12e Conférence internationale sur les représentations d'apprentissage (ICLR 2024) (Vienne, Autriche, 2023) ; <https://openreview.net/forum?id=Ebt7JgMHv1>.
- 1018 D. Stander, Q. Yu, H. Fan, S. Biderman, « Grokking Group Multiplication with Cosets » dans Quarante et unième conférence internationale sur l'apprentissage automatique (2024) ; <https://openreview.net/forum?id=hcQfTsVnBo>.
- 1019 D. Chanin, J. Wilken-Smith, T. Dulka, H. Bhatnagar, J. Bloom, A comme Absorption : étude du fractionnement des caractéristiques et de l'absorption dans les autoencodeurs clairsemés, arXiv [cs.CL] (2024) ; <http://arxiv.org/abs/2409.14507>.
- 1020 J. Adebayo, J. Gilmer, M. Muelly, I. Goodfellow, M. Hardt, B. Kim, « Vérifications de cohérence pour les cartes de saillance » dans *Advances in Neural Information Processing Systems (NeurIPS 2018)* (Curran Associates, Inc., 2018) vol. 31 ; https://proceedings.neurips.cc/paper_files/paper/2018/hash/294a8ed24b1ad22ec2e7efea049b8737-Résumé.html.
- 1021 J. Adebayo, M. Muelly, I. Llicardi, B. Kim, « Tests de débogage pour les explications de modèles » dans *Advances in Neural Information Processing Systems (NeurIPS 2020)* (Curran Associates, Inc., 2020) vol. 33, pp. 700–712 ; <https://proceedings.neurips.cc/paper/2020/hash/075b051ec3d22dac7b33f788da631fd4-Abstract.html>.
- 1022 S. Casper, T. Bu, Y. Li, J. Li, K. Zhang, K. Hariharan, D. Hadfield-Menell, « Red Teaming Deep Neural Networks avec des outils de synthèse de fonctionnalités » dans 37e Conférence sur les systèmes de traitement de l'information neuronale (NeurIPS 2023) (La Nouvelle-Orléans, LA, États-Unis, 2023) ; <https://openreview.net/forum?id=Od6CHhPM7l>.
- 1023 P. Hase, M. Bansal, B. Kim, A. Ghandeharioun, « La localisation informe-t-elle l'édition ? Différences surprenantes entre la localisation basée sur la causalité et l'édition des connaissances dans les modèles linguistiques » dans 37e Conférence sur les systèmes de traitement de l'information neuronale (NeurIPS 2023) (2023) ; <https://openreview.net/forum?id=EldbUIZtbd>.
- 1024 J. Miller, B. Chughtai, W. Saunders, Les mesures de fidélité des circuits de transformateur ne sont pas robustes, arXiv [cs.LG] (2024) ; <http://arxiv.org/abs/2407.08734>.
- 1025* ML Leavitt, A. Morcos, Vers une recherche sur l'interprétabilité falsifiable, arXiv [cs.CY] (2020) ; <http://arxiv.org/abs/2010.12016>.
- 1026* E. Durmus, A. Tamkin, J. Clark, J. Wei, J. Marcus, J. Batson, K. Handa, L. Lovitt, M. Tong, M. McCain, O. Rausch, S. Huang, S. Bowman, S. Ritchie, T. Henighan, D. Ganguli, « Évaluation de l'orientation des fonctionnalités : une étude de cas sur l'atténuation des biais sociaux » (Anthropic, 2024) ; <https://www.anthropic.com/research/evaluating-feature-steering>.
- 1027 GE Hinton, « Représentations distribuées » (CMU-CS-84–157, Université Carnegie-Mellon, 1984) ; <http://shelf2.library.cmu.edu/Tech/19334156.pdf>.
- 1028 Y. Bengio, A. Courville, P. Vincent, Representation Learning: A Review and New Perspectives. *IEEE Transactions on Analyse de modèles et intelligence artificielle* 35, 1798–1828 (2013) ; <https://doi.org/10.1109/TPAMI.2013.50>.
- 1029 L. Gao, J. Schulman, J. Hilton, « Lois d'échelle pour la suroptimisation du modèle de récompense » dans Actes de la 40e Conférence internationale sur l'apprentissage automatique (PMLR, Honolulu, Hawaï, États-Unis, 2023), pp. 10835–10866 ; <https://proceedings.mlr.press/v202/gao23h.html>.
- 1030 P. Singhal, T. Goyal, J. Xu, G. Durrett, Un long chemin à parcourir : étude des corrélations de longueur dans RLHF, arXiv [cs.CL] (2023) ; <http://arxiv.org/abs/2310.03716>.

- 1031 JMV Skalse, NHR Howe, D. Krasheninnikov, D. Krueger, « Définition et caractérisation du jeu de récompense » dans la 36e Conférence sur les systèmes de traitement de l'information neuronale (NeurIPS 2022) (virtuelle, 2022) ; <https://openreview.net/forum?id=yb3HOXO3IX2>.
- 1032 LE McKinney, Y. Duan, D. Krueger, A. Gleave, « Sur la fragilité des fonctions de récompense apprises » dans 36e conférence sur les systèmes de traitement de l'information neuronale (NeurIPS 2022) Atelier sur l'apprentissage par renforcement profond (virtuel, 2022) ; <https://openreview.net/forum?id=9gj9vXfeS-y>.
- 1033 J. Tien, JZ-Y. He, Z. Erickson, A. Dragan, DS Brown, « Confusion causale et mauvaise identification de la récompense dans « Apprentissage par récompense basé sur les préférences » lors de la 11e Conférence internationale sur les représentations d'apprentissage (ICLR 2023) (Kigali, Rwanda, 2022) ; https://openreview.net/forum?id=R0Xxvr_X3ZA.
- 1034 ZX Yong, C. Menghini, S. Bach, « Les langues à faibles ressources jailbreakent GPT-4 » dans l'atelier NeurIPS sur la recherche sur la modélisation linguistique socialement responsable (SoLaR) (La Nouvelle-Orléans, LA, États-Unis, 2023) ; <https://openreview.net/forum?id=pn83r8V2sv>.
- 1035 Y. Huang, L. Sun, H. Wang, S. Wu, Q. Zhang, Y. Li, C. Gao, Y. Huang, W. Lyu, Y. Zhang, X. Li, H. Sun, Z. Liu, Y. Liu, Y. Wang, Z. Zhang, B. Vidgen, ... Y. Zhao, « Position : TrustLLM : la fiabilité dans les grands modèles linguistiques » dans Conférence internationale sur l'apprentissage automatique (PMLR, 2024), pp. 20166–20270 ; <https://proceedings.mlr.press/v235/huang24x.html>.
- 1036 S. Longpre, S. Kapoor, K. Klyman, A. Ramaswami, R. Bommasani, B. Bliili-Hamelin, Y. Huang, A. Skowron, Z.-X. Yong, S. Kotha, Y. Zeng, W. Shi, X. Yang, R. Southen, A. Robey, P. Chao, D. Yang, ... P. Henderson, Une sphère de sécurité pour l'évaluation de l'IA et le Red Teaming, arXiv [cs.AI] (2024) ; <http://arxiv.org/abs/2403.04893>.
- 1037 YM Pa Pa, S. Tanizaki, T. Kou, M. van Eeten, K. Yoshioka, T. Matsumoto, « Le rêve d'un attaquant ? Exploration de la « Capacités de ChatGPT pour le développement de logiciels malveillants » dans les actes du 16e atelier d'expérimentation et de test de cybersécurité (CSET '23) (Association for Computing Machinery, New York, NY, États-Unis, 2023), pp. 10–18 ; <https://doi.org/10.1145/3607505.3607513>.
- 1038 A. Liu, Q. Sheng, X. Hu, « Prévention et détection de la désinformation générée par les grands modèles linguistiques » dans Actes de la 47e Conférence internationale ACM SIGIR sur la recherche et le développement en recherche d'informations (ACM, New York, NY, États-Unis, 2024), pp. 3001–3004 ; <https://doi.org/10.1145/3626772.3661377>.
- 1039 JB Sandbrink, Intelligence artificielle et mauvaise utilisation biologique : différencier les risques des modèles linguistiques et des outils de conception biologique, arXiv [cs.CY] (2023) ; <http://arxiv.org/abs/2306.13952>.
- 1040 L. Pöhler, V. Schrader, A. Ladwein, F. von Keller, Une perspective technologique sur l'utilisation abusive de l'IA disponible, arXiv [cs.CY] (2024) ; <http://arxiv.org/abs/2403.15325>.
- 1041 M. Anderjüng, J. Hazell, Protéger la société contre les abus de l'IA : quand les restrictions sur les capacités sont-elles justifiées ?, arXiv [cs.AI] (2023) ; <http://arxiv.org/abs/2303.09377>.
- 1042 A. Karamolegkou, J. Li, L. Zhou, A. Søgaard, « Violations du droit d'auteur et grands modèles linguistiques » dans Actes de la Conférence 2023 sur les méthodes empiriques en traitement du langage naturel (EMNLP 2023), H. Bouamor, J. Pino, K. Bali, éd. (Association for Computational Linguistics, Singapour, 2023), pp. 7403–7412 ; <https://doi.org/10.18653/v1/2023.emnlp-main.458>.
- 1043 H. Li, D. Guo, W. Fan, M. Xu, J. Huang, F. Meng, Y. Song, « Attaques de confidentialité par jailbreaking en plusieurs étapes sur ChatGPT » dans la conférence 2023 sur les méthodes empiriques en traitement du langage naturel (EMNLP 2023) (Singapour, 2023) ; <https://openreview.net/forum?id=ls4Pfst2jZ>.
- 1044* M. Nasr, N. Carlini, J. Hayase, M. Jagielski, A. Feder Cooper, D. Ippolito, CA Choquette-Choo, E. Wallace, F. Tramèr, K. Lee, Extraction évolutive de données de formation à partir de modèles de langage (de production), arXiv [cs.LG] (2023) ; <http://arxiv.org/abs/2311.17035>.
- 1045 av. J.-C. Das, MH Amini, Y. Wu, Défis de sécurité et de confidentialité des grands modèles linguistiques : une étude, arXiv [cs.CL] (2024) ; <http://arxiv.org/abs/2402.00888>.
- 1046 B. Yan, K. Li, M. Xu, Y. Dong, Y. Zhang, Z. Ren, X. Cheng, Sur la protection de la confidentialité des données des grands modèles linguistiques (LLM) : un aperçu, arXiv [cs.CR] (2024) ; <http://arxiv.org/abs/2403.05156>.
- 1047 Y. Yao, J. Duan, K. Xu, Y. Cai, Z. Sun, Y. Zhang, Une enquête sur la sécurité et la confidentialité du modèle de langage à grande échelle (LLM) : Le bon, la brute et le truand. High-Confidence Computing 4, 100211 (2024) ; <https://doi.org/10.1016/j.hcc.2024.100211>.
- 1048 A. Deshpande, V. Murahari, T. Rajpurohit, A. Kalyan, K. Narasimhan, « Toxicité dans Chatgpt : analyse des modèles de langage assignés à la personne » dans Conclusions de l'Association pour la linguistique computationnelle : EMNLP 2023, H. Bouamor, J. Pino, K. Bali, éd. (Association for Computational Linguistics, Singapour, 2023), pp. 1236–1270 ; <https://doi.org/10.18653/v1/2023.findings-emnlp.88>.
- 1049 Y. Qu, X. Shen, X. He, M. Backes, S. Zannettou, Y. Zhang, « Diffusion dangereuse : sur la génération d'images dangereuses

- et les mêmes haineux issus de modèles texte-image » dans les actes de la conférence ACM SIGSAC 2023 sur la sécurité informatique et des communications (CCS '23) (Association for Computing Machinery, New York, NY, États-Unis, 2023), pp. 3403–3417 ; <https://doi.org/10.1145/3576915.3616679>.
- 1050 Z. Xu, S. Jain, M. Kankanhalli, L'hallucination est inévitable : une limitation innée des grands modèles de langage, arXiv [cs.CL] (2024) ; <http://arxiv.org/abs/2401.11817>.
- 1051* Z. Bai, P. Wang, T. Xiao, T. He, Z. Han, Z. Zhang, MZ Shou, Hallucination de grands modèles de langage multimodaux : une Enquête, arXiv [cs.CV] (2024) ; <http://arxiv.org/abs/2404.18930>.
- 1052 Y. Liu, G. Deng, Z. Xu, Y. Li, Y. Zheng, Y. Zhang, L. Zhao, T. Zhang, K. Wang, Y. Liu, Jailbreaking ChatGPT via Prompt Engineering : une étude empirique, arXiv [cs.SE] (2023) ; <http://arxiv.org/abs/2305.13860>.
- 1053 R. Shah, QF Montixi, S. Pour, A. Tagade, J. Rando, « Jailbreaks de boîte noire évolutifs et transférables pour les modèles de langage via la modulation de personnalité » dans la 37e Conférence sur les systèmes de traitement de l'information neuronale (NeurIPS 2023) Atelier de recherche sur la modélisation linguistique socialement responsable (SoLaR) (La Nouvelle-Orléans, LA, États-Unis, 2023) ; <https://openreview.net/forum?id=x3Ltzq1UFg>.
- 1054 N. Carlini, M. Nasr, CA Choquette-Choo, M. Jagielski, I. Gao, PW Koh, D. Ippolito, F. Tramèr, L. Schmidt, « Are « Les réseaux neuronaux alignés sont-ils alignés de manière antagoniste ? » dans la 37e Conférence sur les systèmes de traitement de l'information neuronale (NeurIPS 2023) (La Nouvelle-Orléans, LA, États-Unis, 2023) ; <https://openreview.net/forum?id=OQqoD8Vc3B>.
- 1055 X. Shen, Z. Chen, M. Backes, Y. Shen, Y. Zhang, « Do Anything Now » : caractérisation et évaluation des invites de jailbreak dans la nature sur de grands modèles de langage, arXiv [cs.CR] (2023) ; <http://arxiv.org/abs/2308.03825>.
- 1056* N. Li, Z. Han, I. Steneker, W. Primack, R. Goodside, H. Zhang, Z. Wang, C. Menghini, S. Yue, Les défenses de LLM ne sont pas Robuste aux jailbreaks humains multi-tours pour le moment, arXiv [cs.LG] (2024) ; <http://arxiv.org/abs/2408.15221>.
- 1057 L. Jiang, K. Rao, S. Han, A. Ettinger, F. Brahman, S. Kumar, N. Mireshghallah, X. Lu, M. Sap, Y. Choi, N. Dziri, « WildTeaming à grande échelle : des jailbreaks dans la nature aux modèles de langage (adversarialement) plus sûrs » dans 38e conférence annuelle sur les systèmes de traitement de l'information neuronale (NeurIPS 2024) (2024) ; <https://openreview.net/pdf?id=n5R6TvBVcX>.
- 1058 Z. Dong, Z. Zhou, C. Yang, J. Shao, Y. Qiao, Attaques, défenses et évaluations pour la sécurité des conversations LLM : une enquête (Association for Computational Linguistics, 2024) ; <https://doi.org/10.18653/v1/2024.naacl-long.375>.
- 1059 M. Andriushchenko, F. Croce, N. Flammarion, Jailbreaking de LLMs de premier plan axés sur la sécurité avec une approche adaptative simple Attaques, arXiv [cs.CR] (2024) ; <http://arxiv.org/abs/2404.02151>.
- 1060 Y. Zeng, H. Lin, J. Zhang, D. Yang, R. Jia, W. Shi, « Comment Johnny peut persuader les LLM de les jailbreaker : repenser la persuasion pour remettre en question la sécurité de l'IA en humanisant les LLM » dans Actes de la 62e réunion annuelle de l'Association for Computational Linguistics (Volume 1 : Long Papers) (Association for Computational Linguistics, Stroudsburg, PA, États-Unis, 2024), pp. 14322–14350 ; <https://doi.org/10.18653/v1/2024.acl-long.773>.
- 1061 AG Chowdhury, MM Islam, V. Kumar, FH Shezan, V. Kumar, V. Jain, A. Chadha, Briser les défenses : Une Étude comparative des attaques sur les grands modèles de langage, arXiv [cs.CR] (2024) ; <http://arxiv.org/abs/2403.04786>.
- 1062 MKB Doumbouya, A. Nandi, G. Poesia, D. Ghilardi, A. Goldie, F. Bianchi, D. Jurafsky, CD Manning, H4rm3l : Une référence dynamique des attaques de jailbreak composables pour l'évaluation de la sécurité LLM, arXiv [cs.CR] (2024) ; <http://arxiv.org/abs/2408.04811>.
- 1063* BRY Huang, M. Li, L. Tang, Jailbreaks sans fin avec l'apprentissage bijectif, arXiv [cs.CL] (2024) ; <http://arxiv.org/abs/2410.01294>.
- 1064 X. Qi, Y. Zeng, T. Xie, P.-Y. Chen, R. Jia, P. Mittal, P. Henderson, « Ajustement précis des modèles de langage alignés « Compromis la sécurité, même lorsque les utilisateurs n'en ont pas l'intention ! » dans la 12e Conférence internationale sur les représentations d'apprentissage (ICLR 2024) (Vienne, Autriche, 2023) ; <https://openreview.net/forum?id=hTEGyKf0dZ>.
- 1065 Q. Zhan, R. Fang, R. Bindu, A. Gupta, T. Hashimoto, D. Kang, « Suppression des protections RLHF dans GPT-4 via un réglage fin » lors de la conférence annuelle 2024 du chapitre nord-américain de l'Association for Computational Linguistics (Mexico, Mexique, 2024) ; <https://doi.org/10.48550/arXiv.2311.05553>.
- 1066 S. Jain, R. Kirk, ES Lubana, RP Dick, H. Tanaka, E. Grefenstette, T. Rocktäschel, DS Krueger, Analyse mécaniste des effets du réglage fin sur des tâches définies de manière procédurale, arXiv [cs.LG] (2023) ; <http://arxiv.org/abs/2311.12786>.
- 1067 X. Yang, X. Wang, Q. Zhang, L. Petzold, WY Wang, X. Zhao, D. Lin, Shadow Alignment : la facilité de subvertir les modèles de langage alignés de manière sûre, arXiv [cs.CL] (2023) ; <http://arxiv.org/abs/2310.02949>.
- 1068 R. Bhardwaj, S. Poria, Désalignement du modèle linguistique : Red-Teaming paramétrique pour exposer les préjugés cachés et Biais, arXiv [cs.CL] (2023) ; <http://arxiv.org/abs/2310.14303>.
- 1069 J. Ji, K. Wang, T. Qiu, B. Chen, J. Zhou, C. Li, H. Lou, Y. Yang, Les modèles de langage résistent à l'alignement, arXiv [cs.CL]

- (2024) ; <http://arxiv.org/abs/2406.06144>.
- 1070 X. Qi, A. Panda, K. Lyu, X. Ma, S. Roy, A. Beirami, P. Mittal, P. Henderson, L'alignement de sécurité devrait être approfondi sur plus de quelques jetons, arXiv [cs.CR] (2024) ; <http://arxiv.org/abs/2406.05946>.
- 1071 S. Hu, Y. Fu, ZS Wu, V. Smith, Régénérer la mémoire des LLM non appris grâce à des attaques de réapprentissage ciblées, arXiv [cs.LG] (2024) ; <http://arxiv.org/abs/2406.13356>.
- 1072 D. Halawi, A. Wei, E. Wallace, TT Wang, N. Haghtalab, J. Steinhardt, « Covert Malicious Finetuning : Challenges in « Sauvegarder l'adaptation du LLM » dans la Conférence internationale sur l'apprentissage automatique (PMLR, 2024), pp. 17298–17312 ; <https://proceedings.mlr.press/v235/halawi24a.html>.
- 1073 R. Greenblatt, F. Roger, D. Krasheninnikov, D. Krueger, « Test de résistance à l'élicitation des capacités avec mot de passe « Modèles verrouillés » dans la 38e conférence annuelle sur les systèmes de traitement de l'information neuronale (NeurIPS 2024) (2024) ; <https://openreview.net/pdf?id=zzOOqD6R1b>.
- 1074 M. Lo, F. Barez, S. Cohen, Les grands modèles linguistiques réapprennent les concepts supprimés (Association for Computational Linguistics, 2024) ; <https://doi.org/10.18653/v1/2024.findings-acl.492>.
- 1075 S. Peng, P.-Y. Chen, MD Hull, DH Chau, « Naviguer dans le paysage de la sécurité : mesurer les risques lors du réglage fin des grands modèles linguistiques » dans 38e conférence annuelle sur les systèmes de traitement de l'information neuronale (NeurIPS 2024) (2024) ; <https://openreview.net/pdf?id=GZnsqBwHAG>.
- 1076 A. Sheshadri, A. Ewart, P. Guo, A. Lynch, C. Wu, V. Hebban, H. Sleight, A. C. Stickland, E. Pérez, D. Hadfield-Menell, S. Casper, La formation contradictoire latente améliore la robustesse aux comportements nuisibles persistants dans les LLM, arXiv [cs.LG] (2024) ; <http://arxiv.org/abs/2407.15549>.
- 1077 S. Xhonneux, A. Sordoni, S. Günemann, G. Gidel, L. Schwinn, « Formation contradictoire efficace dans les LLM avec « Attaques continues » dans la 38e conférence annuelle sur les systèmes de traitement de l'information neuronale (NeurIPS 2024) (2024) ; <https://openreview.net/pdf?id=8jB6sGqvgQ>.
- 1078 L. Schwinn, S. Geisler, Revisiter l'alignement robuste des disjoncteurs, arXiv [cs.CR] (2024) ; <http://arxiv.org/abs/2407.15902>.
- 1079 T. Huang, S. Hu, F. Ilhan, SF Tekin, L. Liu, Attaques et défenses de réglage fin nuisibles pour les grands modèles de langage : Une enquête, arXiv [cs.CR] (2024) ; <http://arxiv.org/abs/2409.18169>.
- 1080 J. Łucki, B. Wei, Y. Huang, P. Henderson, F. Tramèr, J. Rando, Une perspective contradictoire sur le désapprentissage des machines pour la sécurité de l'IA, arXiv [cs.LG] (2024) ; <http://arxiv.org/abs/2409.18025>.
- 1081 Y. Wolf, N. Wies, O. Avnery, Y. Levine, A. Shashua, Limitations fondamentales de l'alignement dans le langage de grande taille Modèles, arXiv [cs.CL] (2023) ; <http://arxiv.org/abs/2304.11082>.
- 1082 T. Tseng, E. McLean, K. Pelrine, TT Wang, A. Gleave, Les IA Go peuvent-elles être robustes face aux adversaires ?, arXiv [cs.LG] (2024) ; <http://arxiv.org/abs/2406.12843>.
- 1083 M. Andriushchenko, N. Flammarion, La formation au refus dans les LLM se généralise-t-elle au passé composé ?, arXiv [cs.CL] (2024) ; <http://arxiv.org/abs/2407.11969>.
- 1084 ID Raji, E. Denton, EM Bender, A. Hanna, A. Paullada, « L'IA et le Tout dans le Monde Entier « Benchmark » dans la 35e Conférence sur les systèmes de traitement de l'information neuronale (NeurIPS 2021) Parcours Ensembles de données et Benchmarks (2e tour) (Virtuel, 2021) ; <https://openreview.net/forum?id=j6NxpQbREA1>.
- 1085 B. Hutchinson, N. Rostamzadeh, C. Greer, K. Heller, V. Prabhakaran, « Lacunes d'évaluation dans l'apprentissage automatique « Pratique » dans les actes de la conférence 2022 de l'ACM sur l'équité, la responsabilité et la transparence (FAccT '22) (Association for Computing Machinery, New York, NY, États-Unis, 2022), pp. 1859–1876 ; <https://doi.org/10.1145/3531146.3533233>.
- 1086 S. Casper, C. Ezell, C. Siegmann, N. Kolt, TL Curtis, B. Bucknall, A. Haupt, K. Wei, J. Scheurer, M. Hobbhahn, L. Sharkey, S. Krishna, M. Von Hagen, S. Alberti, A. Chan, Q. Sun, M. Gerovitch, ... D. Hadfield-Menell, L'accès à la boîte noire est insuffisant pour des audits d'IA rigoureux, arXiv [cs.CY] (2024) ; <http://arxiv.org/abs/2401.14446>.
- 1087 B. Ram, P. Verma, Étude de Chatbot basée sur l'intelligence artificielle de ChatGPT, Google AI Bard et Baidu AI. Monde Journal des sciences et technologies avancées de l'ingénierie 8, 258–261 (2023) ; <https://doi.org/10.30574/wjaets.2023.8.1.0045>.
- 1088 MM Maas, « La gouvernance de l'intelligence artificielle en mutation : fondements, facettes, cadres », thèse, Université de Copenhague (2020) ; [https://matthijsmaas.com/uploads/Maas%20-%202021%20-%20Thèse de doctorat - La gouvernance de l'intelligence artificielle en pleine mutation %20monographie.pdf](https://matthijsmaas.com/uploads/Maas%20-%202021%20-%20Thèse%20de%20doctorat%20-La%20gouvernance%20de%20l'intelligence%20artificielle%20en%20pleine%20mutation%20monographie.pdf).
- 1089 PM Napoli, Les médias sociaux et l'intérêt public (Columbia University Press, 2019) ; <https://cup.columbia.edu/book/social-media-and-the-public-interest/9780231184540>.
- 1090 JM Balkin, Comment réglementer (et ne pas réglementer) les médias sociaux. Journal of Free Speech Law 1, 71–96 (2021) ;

- <https://doi.org/10.2139/ssrn.3484114>.
- 1091 RH Frank, PJ Cook, Marchés où le gagnant rafle tout. *Studies in Microeconomics* 1, 131–154 (2013) ; <https://doi.org/10.1177/2321022213501254>.
- 1092 BA Prakash, A. Beutel, R. Rosenfeld, C. Faloutsos, « Le gagnant rafle tout : virus concurrents ou idées sur le fair-play Réseaux » dans les actes de la 21^e Conférence internationale sur le World Wide Web - WWW '12 (ACM Press, New York, New York, États-Unis, 2012) ; <https://doi.org/10.1145/2187836.2187975>.
- 1093 TA Han, LM Pereira, T. Lenaerts, « Modélisation et influence de la guerre des enchères sur l'IA : un programme de recherche » dans Actes de la conférence 2019 de l'AAAI/ACM sur l'IA, l'éthique et la société (AIES '19) (New York, NY, États-Unis, 2019), pp. 5–11 ; <https://doi.org/10.1145/3306618.3314265>.
- 1094 T. Cimpeanu, FC Santos, LM Pereira, T. Lenaerts, TA Han, Courses de développement de l'intelligence artificielle en Paramètres hétérogènes. *Scientific Reports* 12, 1723 (2022) ; <https://doi.org/10.1038/s41598-022-05729-3>.
- 1095 A. Guasti, M. Koenig-Archibugi, La concurrence commerciale mondiale a-t-elle vraiment conduit à une course vers le bas du marché du travail ? Normes ? *International Studies Quarterly* : une publication de l'International Studies Association 66, sqac061 (2022) ; <https://doi.org/10.1093/isq/sqac061>.
- 1096 G. Porter, Concurrence commerciale et normes de pollution : « course vers le bas » ou « coincé au bas ». *Journal of Environment & Development* 8, 133–151 (1999) ; <https://doi.org/10.1177/107049659900800203>.
- 1097 D. Vera, C. Rusche, « L'économie des plateformes » (Institut der deutschen Wirtschaft, 2018) ; <https://www.iwkoeln.de/en/studies/vera-demary-christian-rusche-the-Economics-of-platforms.html>.
- 1098 MF Niculescu, DJ Wu, L. Xu, Partage stratégique de la propriété intellectuelle : concurrence sur une plateforme technologique ouverte sous effets de réseau. *Information Systems Research* : ISR 29, 498–519 (2018) ; <https://doi.org/10.1287/isre.2017.0756>.
- 1099 NL Rose, La peur de l'avion ? Analyses économiques de la sécurité aérienne. *The Journal of Economic Perspectives* : A Journal of the American Economic Association 6, 75–94 (1992) ; <https://doi.org/10.1257/jep.6.2.75>.
- 1100 J. Tirole, La théorie de l'organisation industrielle (MIT Press, Londres, Angleterre, 1988).
- 1101 S. Armstrong, N. Bostrom, C. Shulman, La course au précipice : un modèle de développement de l'intelligence artificielle. *AI & Society* 31, 201–206 (2016) ; <https://doi.org/10.1007/s00146-015-0590-y>.
- 1102 GH Stern, RJ Feldman, Trop gros pour faire faillite : les dangers des renfouements bancaires (Brookings Institution Press, 2009) ; <https://www.brookings.edu/books/too-big-to-fail/>.
- 1103 BE Gup, Financial Management Association International, Too Big to Fail : Politiques et pratiques du gouvernement Plans de sauvetage (Praeger, Westport, Conn, éd. 1, 2003) ; https://library-search.open.ac.uk/permalink/44OPN_INST/la9sg5/alma9952597297902316.
- 1104 V. Acharya, D. Anginer, JA Warburton, « La fin de la discipline du marché ? Les attentes des investisseurs en matière de discipline implicite « Garanties gouvernementales » (2022) ; <https://cepr.org/publications/dp17426>.
- 1105 K. Pernell, J. Jung, Repenser le risque moral : protection gouvernementale et prise de risque bancaire. *Socio-Economic Review* 22, 625–653 (2024) ; <https://doi.org/10.1093/ser/mwad050>.
- 1106 WJ Baumol, WE Oates, La théorie de la politique environnementale (Cambridge University Press, Cambridge, Angleterre, éd. 2, 1988) ; <https://doi.org/10.1017/cbo9781139173513>.
- 1107 P. DeCicca, D. Kenkel, MF Lovenheim, L'économie de la réglementation du tabac : une analyse approfondie. *Journal of Economic Literature* 60, 883–970 (2022) ; <https://doi.org/10.1257/jel.20201482>.
- 1108 J. Guerreiro, S. Rebelo, P. Teles, « Régulation de l'intelligence artificielle » (w31921, Bureau national de recherche économique, 2023) ; <https://doi.org/10.3386/w31921>.
- 1109 L. Dallas, « Le court-termisme, la crise financière et la gouvernance d'entreprise » (Faculté de droit de l'Université de San Diego, 2012) ; <http://dx.doi.org/>.
- 1110 N. Kolt, M. Anderjunga, J. Barnhart, A. Brass, K. Esvelt, GK Hadfield, L. Heim, M. Rodriguez, JB Sandbrink, T. Woodside, Reporting responsable pour le développement de l'IA de pointe. *arXiv [cs.CY]* (2024) ; <http://arxiv.org/abs/2404.02675>.
- 1111 M. Anderjunga, J. Barnhart, A. Korinek, J. Leung, C. O'Keefe, J. Whittlestone, S. Avin, M. Brundage, J. Bullock, D. Cass-Beggs, B. Chang, T. Collins, T. Fist, G. Hadfield, A. Hayes, L. Ho, S. Hooker, ... K. Wolf, Frontier AI Regulation: Managing Emerging Risks to Public Safety, *arXiv [cs.CY]* (2023) ; <http://arxiv.org/abs/2307.03718>.
- 1112 L. Collina, M. Sayyadi, M. Provitera, Réponses aux questions critiques sur la responsabilité de l'IA. *California Management Review Insights* (2023) ; <https://cmr.berkeley.edu/2023/11/critical-issues-about-ai-accountability-answered/>.
- 1113 AT da Fonseca, E. Vaz de Sequeira, L. Barreto Xavier, « Responsabilité des systèmes pilotés par l'IA » dans Perspectives multidisciplinaires sur l'intelligence artificielle et le droit, H. Sousa Antunes, PM Freitas, AL Oliveira, C. Martins

- Pereira, E. Vaz de Sequeira, L. Barreto Xavier, éd. (Springer International Publishing, Cham, 2024), pp. https://doi.org/10.1007/978-3-031-41264-6_16.
- 1114 M. Buiten, A. de Streel, M. Peitz, Le droit et l'économie de la responsabilité de l'IA. Rapport sur le droit et la sécurité informatique 48, 105794 (2023) ; <https://doi.org/10.1016/j.clsr.2023.105794>.
- 1115 T. Miller, Explication en intelligence artificielle : perspectives des sciences sociales. *Intelligence artificielle* 267, 1–38 (2019) ; <https://doi.org/10.1016/j.artint.2018.07.007>.
- 1116 F. Doshi-Velez, B. Kim, Vers une science rigoureuse de l'apprentissage automatique interprétable, *arXiv [stat.ML]* (2017) ; <http://arxiv.org/abs/1702.08608>.
- 1117 ZC Lipton, Le mythe de l'interprétabilité des modèles : dans l'apprentissage automatique, le concept d'interprétabilité est à la fois important et délicat. *ACM Queue : Tomorrow's Computing Today* 16, 31–57 (2018) ; <https://doi.org/10.1145/3236386.3241340>.
- 1118 T. Räuker, A. Ho, S. Casper, D. Hadfield-Menell, Vers une IA transparente : une enquête sur l'interprétation des structures internes des réseaux neuronaux profonds, *arXiv [cs.LG]* (2022) ; <http://arxiv.org/abs/2207.13243>.
- 1119 M. Busuioc, Intelligence artificielle responsable : responsabiliser les algorithmes. *Public Administration Review* 81, 825–836 (2021) ; <https://doi.org/10.1111/puar.13293>.
- 1120 F. Doshi-Velez, M. Kortz, R. Budish, C. Bavitz, S.J. Gershman, D. O'Brien, K. Scott, S. Shieber, J. Waldo, D. Weinberger, A. Weller, A. Wood, « Responsabilité de l'IA en vertu de la loi : le rôle de l'explication » (Groupe de travail du Centre Berkman Klein sur l'explication et la loi, 2017) ; <http://nrs.harvard.edu/urn-3:HUL.InstRepos:34372584>.
- 1121 R. Palin, I. Habli, « Assurance de la sécurité automobile – Une approche fondée sur des cas de sécurité » dans *Computer Safety, Reliability, and Security (SAFECOMP 2010)*, E. Schoitsch, éd. (Springer, Berlin, Heidelberg, 2010) Notes de cours en informatique (LNPS), pp. 82–96 ; https://doi.org/10.1007/978-3-642-15651-9_7.
- 1122 I. Livshitz, PA Lontsikh, NP Lontsikh, EY Golovina, OM Safonova, « Une étude des méthodes modernes de gestion des risques pour l'assurance de la sécurité industrielle dans l'industrie des carburants et de l'énergie » dans la Conférence internationale 2021 sur la gestion de la qualité, la sécurité des transports et de l'information, les technologies de l'information (IT&QM&IS) (2021), pp. 165–167 ; <https://doi.org/10.1109/ITQMIS53292.2021.9642791>.
- 1123 M.L. Cummings, Repenser la maturité de l'intelligence artificielle dans les environnements critiques pour la sécurité. *AI Magazine* 42, 6–15 (2021) ; <https://ojs.aaai.org/aimagazine/index.php/aimagazine/article/view/7394>.
- 1124 N. Kolt, Gouverner les agents d'IA (2024) ; <https://doi.org/10.2139/ssrn.4772956>.
- 1125 P. Verdegem, Démanteler le capitalisme de l'IA : les biens communs comme alternative à la concentration du pouvoir des grandes entreprises technologiques. *AI & Society* 39, 1–11 (2022) ; <https://doi.org/10.1007/s00146-022-01437-8>.
- 1126 K. Crawford, Atlas de l'IA, Yale University Press Londres (2021) ; <https://yalebooks.co.uk/9780300264630/atlas-of-ai>.
- 1127 J. Angwin, A. Nelson, R. Palta, « Vous recherchez des informations électorales fiables ? Ne faites pas confiance à l'IA » (The AI Democracy Projects, 2024) ; <https://www.proofnews.org/seeking-election-information-dont-trust-ai/>.
- 1128 H. Shen, A. DeVos, M. Eslami, K. Holstein, Audit d'algorithmes au quotidien : comprendre le pouvoir de l'audit quotidien Utilisateurs dans la détection de comportements algorithmiques nuisibles. *Actes de l'ACM sur l'interaction homme-machine* 5, 1–29 (2021) ; <https://doi.org/10.1145/3479577>.
- 1129 G. Abercrombie, D. Benbouzid, P. Giudici, D. Golpayegani, J. Hernandez, P. Noro, H. Pandit, E. Paraschou, C. Pownall, J. Prajapati, MA Sayre, U. Sengupta, A. Suriyawongkul, R. Thelot, S. Vei, L. Waltersdorfer, Une taxonomie collaborative centrée sur l'humain des dommages causés par l'IA, les algorithmes et l'automatisation, *arXiv [cs.LG]* (2024) ; <http://arxiv.org/abs/2407.01294>.
- 1130 J. Molloy, S. Shahbeigi, JA McDermid, Analyse des risques et de la sécurité de la perception basée sur l'apprentissage automatique Capacités dans les véhicules autonomes. *Computer* 57, 60–70 (2024) ; <https://doi.org/10.1109/mc.2024.3443751>.
- 1131 Y. Jia, T. Lawton, J. Burden, J. McDermid, I. Habli, Conception axée sur la sécurité de l'apprentissage automatique pour le traitement du sepsis. *Journal d'informatique biomédicale* 117, 103762 (2021) ; <https://doi.org/10.1016/j.jbi.2021.103762>.
- 1132 R. Hawkins, C. Picardi, L. Donnell, M. Ireland, Création d'un dossier d'assurance de sécurité pour un système de détection et d'alerte des incendies de forêt par satellite basé sur l'apprentissage automatique. *Journal of Intelligent & Robotic Systems* 108, 1–21 (2023) ; <https://doi.org/10.1007/s10846-023-01905-3>.
- 1133 P. Festor, Y. Jia, AC Gordon, AA Faisal, I. Habli, M. Komorowski, Assurer la sécurité des tests cliniques basés sur l'IA Systèmes d'aide à la décision : une étude de cas du clinicien IA pour le traitement du sepsis. *BMJ Health & Care Informatics* 29, e100549 (2022) ; <https://doi.org/10.1136/bmjhci-2022-100549>.
- 1134 Département des sciences, de l'innovation et de la technologie, « Engagements en matière de sécurité de l'IA de pointe, AI Seoul Summit 2024 » (GOV.Royaume-Uni, 2024) ; <https://www.gov.uk/government/publications/frontier-ai-safety-commitments-ai-seoul->

- sommet-2024/engagements-en-matière-de-sécurité-de-l'IA-frontière-ai-seoul-sommet-2024.
- 1135 R. Schwartz, J. Fiscus, K. Greene, G. Waters, R. Chowdhury, T. Jensen, C. Greenberg, A. Godil, R. Amironesei, P. Hall, S. Jain, « Le plan d'évaluation pilote du NIST pour l'évaluation des risques et des impacts de l'IA (ARIA) » (Institut national des normes et de la technologie des États-Unis, 2024) ; <https://ai-challenges.nist.gov/uassets/7>.
- 1136 CG Northcutt, A. Athalye, J. Mueller, « Les erreurs d'étiquetage généralisées dans les ensembles de tests déstabilisent l'apprentissage automatique « Benchmarks » dans la 35e Conférence sur les systèmes de traitement de l'information neuronale (NeurIPS 2021) Parcours Ensembles de données et Benchmarks (1er tour) (Virtuel, 2021) ; <https://openreview.net/forum?id=XccDXrDNLeK>.
- 1137 Z. Xiao, S. Zhang, V. Lai, QV Liao, Évaluation des mesures d'évaluation : un cadre pour l'analyse des mesures d'évaluation NLG à l'aide de la théorie de la mesure (Association for Computational Linguistics, 2023) ; <https://doi.org/10.18653/v1/2023.emnlp-main.676>.
- 1138 M. Sclar, Y. Choi, Y. Tsvetkov, A. Suhr, Quantification de la sensibilité des modèles de langage aux caractéristiques parasites dans la conception des invites ou : comment j'ai appris à commencer à m'inquiéter du formatage des invites, arXiv [cs.CL] (2023) ; <http://arxiv.org/abs/2310.11324>.
- 1139 B. Shu, L. Zhang, M. Choi, L. Dunagan, L. Logeswaran, M. Lee, D. Card, D. Jurgens, « Vous n'avez pas besoin d'une personnalité « Test pour savoir si ces modèles ne sont pas fiables : évaluation de la fiabilité des grands modèles linguistiques sur des instruments psychométriques » dans Actes de la conférence 2024 du chapitre nord-américain de l'Association for Computational Linguistics : Human Language Technologies (Volume 1 : Long Papers) (Association for Computational Linguistics, Stroudsburg, PA, États-Unis, 2024), pp. 5263–5281 ; <https://doi.org/10.18653/v1/2024.naacl-long.295>.
- 1140 A. Bavaresco, R. Bernardi, L. Bertolazzi, D. Elliott, R. Fernández, A. Gatt, E. Ghaleb, M. Giulianelli, M. Hanna, A. Koller, AFT Martins, P. Mondorf, V. Neplenbroek, S. Pezzelle, B. Plank, D. Schlangen, A. Suglia, ... A. Testoni, LLM au lieu de Human Des juges ? Une étude empirique à grande échelle portant sur 20 tâches d'évaluation de la PNL, arXiv [cs.CL] (2024) ; <http://arxiv.org/abs/2406.18403>.
- 1141 ISACA, « Le cadre de gestion des risques informatiques » (2009) ; https://www.hci-til.com/ITIL_v3/docs/RiskIT_FW_30June2010_Research.pdf.
- 1142 US AI Safety Institute, UK AI Safety Institute, « Test de pré-déploiement conjoint AISI des États-Unis et AISI du Royaume-Uni » (National Institute of Standards and Technology ; Département de la science, de l'innovation et de la technologie, 2024) ; <https://www.nist.gov/system/files/documents/2024/11/19/Upgraded%20Sonnet-Publication-US.pdf>.
- 1143 G. Leech, JJ Vazquez, N. Kupper, M. Yagudin, L. Aitchison, Pratiques douteuses en apprentissage automatique, arXiv [cs.LG] (2024) ; <http://arxiv.org/abs/2407.12220>.
- 1144* L. Madaan, AK Singh, R. Schaeffer, A. Poulton, S. Koyejo, P. Stenatorp, S. Narang, D. Hupkes, Quantifying Variance dans les repères d'évaluation, arXiv [cs.LG] (2024) ; <http://arxiv.org/abs/2406.10229>.
- 1145 C. Xu, S. Guan, D. Greene, M.-T. Kechadi, Contamination des données de référence des grands modèles linguistiques : une enquête, arXiv [cs.CL] (2024) ; <http://arxiv.org/abs/2406.04244>.
- 1146 Y. Chang, X. Wang, J. Wang, Y. Wu, L. Yang, K. Zhu, H. Chen, X. Yi, C. Wang, Y. Wang, W. Ye, Y. Zhang, Y. Chang, PS Yu, Q. Yang, X. Xie, Une enquête sur l'évaluation des grands modèles linguistiques. ACM Transactions on Intelligent Systems and Technology 15, 39 : 1–39 : 45 (2024) ; <https://doi.org/10.1145/3641289>.
- 1147* W. Zhong, R. Cui, Y. Guo, Y. Liang, S. Lu, Y. Wang, A. Saied, W. Chen, N. Duan, AGIEval : une référence centrée sur l'humain pour évaluer les modèles de fondation, arXiv [cs.CL] (2023) ; <http://arxiv.org/abs/2304.06364>.
- 1148 L. Zheng, W.-L. Chiang, Y. Sheng, S. Zhuang, Z. Wu, Y. Zhuang, Z. Lin, Z. Li, D. Li, E. Xing, H. Zhang, JE Gonzalez, I. Stoica, « Juger LLM-as-a-Judge avec MT-Bench et Chatbot Arena » dans la 37e Conférence sur les systèmes de traitement de l'information neuronale (NeurIPS 2023) Datasets and Benchmarks Track (La Nouvelle-Orléans, LA, États-Unis, 2023) ; <https://openreview.net/forum?id=uccHPGDlao>.
- 1149* S. Yao, N. Shinn, P. Razavi, K. Narasimhan, τ-Bench : une référence pour l'interaction outil-agent-utilisateur dans le monde réel Domaines, arXiv [cs.AI] (2024) ; <http://arxiv.org/abs/2406.12045>.
- 1150 P. Liang, R. Bommasani, T. Lee, D. Tsipras, D. Soylu, M. Yasunaga, Y. Zhang, D. Narayanan, Y. Wu, A. Kumar, B. Newman, B. Yuan, B. Yan, C. Zhang, CA Cosgrove, CD Manning, C. Re, ... Y. Koreeda, Évaluation holistique des modèles linguistiques. Transactions on Machine Learning Research (2023) ; <https://openreview.net/forum?id=iO4LZibEqW>.
- 1151 A. Reuel, A. Hardy, C. Smith, M. Lamparth, M. Hardy, MJ Kochenderfer, BetterBench : évaluation des repères d'IA, découverte des problèmes et établissement des meilleures pratiques, arXiv [cs.AI] (2024) ; <http://arxiv.org/abs/2411.12990>.
- 1152* E. Miller, Ajout de barres d'erreur aux évaluations : une approche statistique des évaluations de modèles linguistiques, arXiv [stat.AP] (2024) ; <http://arxiv.org/abs/2411.00640>.
- 1153* N. Sambasivan, E. Arnesen, B. Hutchinson, V. Prabhakaran, Non-portabilité de l'équité algorithmique en Inde, arXiv

- [cs.CY] (2020) ; <http://arxiv.org/abs/2012.03659>.
- 1154 IO Gallegos, RA Rossi, J. Barrow, MM Tanjim, S. Kim, F. Derroncourt, T. Yu, R. Zhang, NK Ahmed, Biais et équité dans les grands modèles linguistiques : une enquête. *Linguistique computationnelle (Association for Computational Linguistics)* 50, 1–83 (2024) ; https://doi.org/10.1162/coli_a_00524.
- 1155 K. Charmaz, *Construire une théorie ancrée* (SAGE Publications, Thousand Oaks, CA, 2014).
- 1156 T. Shin, Y. Razeghi, RL Logan IV, E. Wallace, S. Singh, « AutoPrompt : extraction de connaissances à partir de modèles linguistiques avec des invites générées automatiquement » dans *Actes de la Conférence 2020 sur les méthodes empiriques en traitement du langage naturel (EMNLP 2020)*, B. Webber, T. Cohn, Y. He, Y. Liu, éd. (Association for Computational Linguistics, en ligne, 2020), pp. 4222–4235 ; <https://doi.org/10.18653/v1/2020.emnlp-main.346>.
- 1157 E. Perez, S. Huang, F. Song, T. Cai, R. Ring, J. Aslanides, A. Glaese, N. McAleese, G. Irving, « Red Teaming Language Models with Language Models » dans *Actes de la Conférence 2022 sur les méthodes empiriques en traitement du langage naturel (EMNLP 2022)*, Y. Goldberg, Z. Kozareva, Y. Zhang, éd. (Association for Computational Linguistics, Abu Dhabi, Émirats arabes unis, 2022), pp. 3419–3448 ; <https://doi.org/10.18653/v1/2022.emnlp-main.225>.
- 1158* D. Ganguli, L. Lovitt, J. Kernion, A. Askell, Y. Bai, S. Kadavath, B. Mann, E. Perez, N. Schiefer, K. Ndousse, A. Jones, S. Bowman, A. Chen, T. Conerly, N. DasSarma, D. Drain, N. Elhage, ... J. Clark, « Modèles linguistiques Red Teaming pour réduire les dommages : méthodes, mise à l'échelle des comportements et leçons apprises » (*Anthropic*, 2022) ; <http://arxiv.org/abs/2209.07858>.
- 1159 S. Casper, J. Lin, J. Kwon, G. Culp, D. Hadfield-Menell, Explorer, établir, exploiter : modèles de langage Red Teaming à partir de zéro, *arXiv [cs.CL]* (2023) ; <http://arxiv.org/abs/2306.09442>.
- 1160 S. Tong, E. Jones, J. Steinhardt, « Échecs de production de masse des systèmes multimodaux avec des modèles linguistiques » dans *37e Conférence sur les systèmes de traitement de l'information neuronale (NeurIPS 2023)* (La Nouvelle-Orléans, LA, États-Unis, 2023) ; <https://openreview.net/forum?id=T6iiOqsGOh>.
- 1161 M. Mazeika, L. Phan, X. Yin, A. Zou, Z. Wang, N. Mu, E. Sakhaee, N. Li, S. Basart, B. Li, D. Forsyth, D. Hendrycks, HarmBench : un cadre d'évaluation standardisé pour le red teaming automatisé et le refus robuste, *arXiv [cs.LG]* (2024) ; <http://arxiv.org/abs/2402.04249>.
- 1162 P. Chao, A. Robey, E. Dobriban, H. Hassani, GJ Pappas, E. Wong, Jailbreaking Black Box Modèles de langage volumineux dans Vingt requêtes, *arXiv [cs.LG]* (2023) ; <http://arxiv.org/abs/2310.08419>.
- 1163 D. Ziegler, S. Nix, L. Chan, T. Bauman, P. Schmidt-Nielsen, T. Lin, A. Scherlis, N. Nabeshima, B. Weinstein-Raun, D. de Haas, B. Shlegeris, N. Thomas, « Formation contradictoire pour une fiabilité à enjeux élevés » dans *Advances in Neural Information Processing Systems (NeurIPS 2022)* (La Nouvelle-Orléans, LA, États-Unis, 2022) vol. 35, pp. 9274–9286 ; https://proceedings.neurips.cc/paper_files/paper/2022/hash/3c44405d619a6920384a45bce876b41e-Abstract-Conference.html.
- 1164 A. Rao, S. Vashista, A. Naik, S. Aditya, M. Choudhury, « Tromper les LLM dans la désobéissance : formalisation, analyse et détection des jailbreaks » dans *Conférence internationale conjointe 2024 sur la linguistique computationnelle, les ressources linguistiques et l'évaluation (LREC-COLING 2024)* (Turin, Italie, 2024) ; <https://doi.org/10.48550/arXiv.2305.14965>.
- 1165* A. Mehrotra, M. Zampetakis, P. Kassianik, B. Nelson, H. Anderson, Y. Singer, A. Karbasi, Tree of Attacks : jailbreak automatique des LLM Black-Box, *arXiv [cs.LG]* (2023) ; <http://arxiv.org/abs/2312.02119>.
- 1166 TD Pala, VYH Toh, R. Bhardwaj, S. Poria, Ferret : Red Teaming automatisé plus rapide et efficace avec une technique de notation basée sur les récompenses, *arXiv [cs.CL]* (2024) ; <http://arxiv.org/abs/2408.10701>.
- 1167 M. Feffer, A. Sinha, ZC Lipton, H. Heidari, Red-Teaming pour l'IA générative : solution miracle ou théâtre de sécurité ?, *arXiv [cs.CY]* (2024) ; <http://arxiv.org/abs/2401.15897>.
- 1168* L. Weidinger, J. Mellor, BG Pegueroles, N. Marchal, R. Kumar, K. Lum, C. Akbulut, M. Diaz, S. Bergman, M. Rodriguez, V. Rieser, W. Isaac, STAR : Approche sociotechnique des modèles de langage Red Teaming, *arXiv [cs.AI]* (2024) ; <http://arxiv.org/abs/2406.11757>.
- 1169 P. Chao, E. DeBenedetti, A. Robey, M. Andriushchenko, F. Croce, V. Sehwag, E. Dobriban, N. Flammarion, GJ Pappas, F. Tramèr, H. Hassani, E. Wong, JailbreakBench : une référence de robustesse ouverte pour le jailbreaking de grands modèles de langage, *arXiv [cs.CR]* (2024) ; <http://arxiv.org/abs/2404.01318>.
- 1170 US AI Safety Institute, « Gestion des risques d'utilisation abusive pour les modèles de fondation à double usage » (NIST, 2024) ; <https://doi.org/10.6028/nist.ai.800-1.ipd>.
- 1171 W. Tann, Y. Liu, JH Sim, CM Seah, E.-C. Chang, Utilisation de grands modèles de langage pour la capture de la cybersécurité Défis liés aux drapeaux et questions de certification, *arXiv [cs.AI]* (2023) ; <http://arxiv.org/abs/2308.10443>.
- 1172 D. Kang, X. Li, I. Stoica, C. Guestrin, M. Zaharia, T. Hashimoto, « Exploiter le comportement programmatique des LLM : double-

- « Utilisation via des attaques de sécurité standard » dans les ateliers IEEE sur la sécurité et la confidentialité (SPW) 2024 (IEEE, 2024), pp. 132–143 ; <https://doi.org/10.1109/spw63631.2024.00018>.
- 1173 FN Motlagh, M. Hajizadeh, M. Majd, P. Najafi, F. Cheng, C. Meinel, Grands modèles de langage en cybersécurité : État de l'art, arXiv [cs.CR] (2024) ; <http://arxiv.org/abs/2402.00891>.
- 1174 A. Hagerty, I. Rubinov, Éthique mondiale de l'IA : un examen des impacts sociaux et des implications éthiques de l'intelligence artificielle, arXiv [cs.CY] (2019) ; <http://arxiv.org/abs/1907.07892>.
- 1175 MM Maas, « Aligner la réglementation de l'IA sur les changements sociotechniques » dans The Oxford Handbook of AI Governance, JB Bullock, Y.-C. Chen, J. Himmelreich, VM Hudson, A. Korinek, MM Young, B. Zhang, éd. (Oxford University Press, 2022) ; <https://doi.org/10.1093/oxfordhb/9780197579329.013.22>.
- 1176 D. Dalrymple, J. Skalse, Y. Bengio, S. Russell, M. Tegmark, S. Seshia, S. Omohundro, C. Szegedy, B. Goldhaber, N. Ammann, A. Abate, J. Halpern, C. Barrett, D. Zhao, T. Zhi-Xuan, J. Wing, J. Tenenbaum, Vers une IA sûre garantie : un cadre pour garantir des systèmes d'IA robustes et fiables, arXiv [cs.AI] (2024) ; <http://arxiv.org/abs/2405.06624>.
- 1177 A. Reuel, B. Bucknall, S. Casper, T. Fist, L. Soder, O. Aarne, L. Hammond, L. Ibrahim, A. Chan, P. Wills, M. Anderljung, B. Garfinkel, L. Heim, A. Trask, G. Mukobi, R. Schaeffer, M. Baker, ... R. Trager, Problèmes ouverts dans la gouvernance technique de l'IA, arXiv [cs.CY] (2024) ; <http://arxiv.org/abs/2407.14981>.
- 1178 R. Ren, S. Basart, A. Khoja, A. Gatti, L. Phan, X. Yin, M. Mazeika, A. Pan, G. Mukobi, RH Kim, S. Fitz, D. Hendrycks, « Safetywashing : les repères de sécurité de l'IA mesurent-ils réellement les progrès en matière de sécurité ? » dans 38e conférence sur les ensembles de données et les repères des systèmes de traitement de l'information neuronale (2024) ; <https://openreview.net/pdf?id=YagfTP3RK6>.
- 1179 BS Bucknall, RF Trager, « Accès structuré pour la recherche tierce sur les modèles d'IA de pointe : étude des exigences d'accès aux modèles des chercheurs » (Oxford Martin School, Université d'Oxford et Centre pour la gouvernance de l'IA, 2023) ; https://cdn.governance.ai/Structured_Access_for_Third-Party_Research.pdf.
- 1180 A. Birhane, VU Prabhu, E. Kahembwe, Ensembles de données multimodaux : misogynie, pornographie et stéréotypes malins, arXiv [cs.CY] (2021) ; <http://arxiv.org/abs/2110.01963>.
- 1181 R. Ashmore, R. Calinescu, C. Paterson, Assurer le cycle de vie de l'apprentissage automatique. ACM Computing Surveys 54, 1–39 (2022) ; <https://doi.org/10.1145/3453444>.
- 1182 S. Casper, X. Davies, C. Shi, TK Gilbert, J. Scheurer, J. Rando, R. Freedman, T. Korbak, D. Lindner, P. Freire, TT Wang, S. Marks, C.-R. Segerie, M. Carroll, A. Peng, P. Christoffersen, M. Damani, ... D. Hadfield-Menell, Problèmes ouverts et limites fondamentales de l'apprentissage par renforcement à partir de rétroactions humaines. Transactions on Machine Learning Research (2023) ; <https://openreview.net/forum?id=bx24KpJ4Eb>.
- 1183 T. Shevlane, Accès structuré : un paradigme émergent pour un déploiement sûr de l'IA, arXiv [cs.AI] (2022) ; <http://arxiv.org/abs/2201.05159>.
- 1184 J. Petrie, O. Aarne, N. Amman, D. Dalrymple, Rapport intermédiaire : Mécanismes pour des garanties flexibles basées sur le matériel. (2024) ; https://yoshuabengio.org/wp-content/uploads/2024/09/FlexHEG-Interim-Report_2024.pdf.
- 1185 S. Costanza-Chock, ID Raji, J. Buolamwini, « Qui vérifie les auditeurs ? Recommandations issues d'une analyse de terrain de l'écosystème de l'audit algorithmique » dans Actes de la conférence 2022 de l'ACM sur l'équité, la responsabilité et la transparence (FAccT '22) (Association for Computing Machinery, New York, NY, États-Unis, 2022), pp. 1571–1583 ; <https://doi.org/10.1145/3531146.3533213>.
- 1186 M. Feffer, M. Skirpan, Z. Lipton, H. Heidari, « De l'élicitation des préférences à l'apprentissage automatique participatif : une enquête critique et des lignes directrices pour la recherche future » dans Actes de la conférence 2023 de l'AAAI/ACM sur l'IA, l'éthique et la société (AIES '23) (ACM, Montréal QC Canada, 2023), pp. 38–48 ; <https://doi.org/10.1145/3600211.3604661>.
- 1187 F. Delgado, S. Yang, M. Madaio, Q. Yang, « Le tournant participatif dans la conception de l'IA : fondements théoriques et état actuel de la pratique » dans Actes de la 3e conférence de l'ACM sur l'équité et l'accès aux algorithmes, mécanismes et optimisation (EAAMO '23) (Association for Computing Machinery, New York, NY, États-Unis, 2023), pp. 1–23 ; <https://doi.org/10.1145/3617694.3623261>.
- 1188 J. Metcalf, E. Moss, EA Watkins, R. Singh, MC Elish, « Évaluations d'impact algorithmique et responsabilité : la « Co-construction des impacts » dans les actes de la conférence 2021 de l'ACM sur l'équité, la responsabilité et la transparence (FAccT '21) (Association for Computing Machinery, New York, NY, États-Unis, 2021), pp. 735–746 ; <https://doi.org/10.1145/3442188.3445935>.
- 1189 D. Martin Jr, V. Prabhakaran, J. Kuhlberg, A. Smart, WS Isaac, « Formulation participative des problèmes pour une économie plus juste » « Machine Learning Through Community Based System Dynamics » dans l'atelier ICLR sur l'apprentissage automatique dans la vie réelle (2020) ; <https://doi.org/10.48550/arXiv.2005.07572>.

- 1190 S. Fazelpour, M. De-Arteaga, Diversité dans les systèmes d'apprentissage automatique sociotechnique. *Big Data & Society* 9, 205395172210820 (2022) ; <https://doi.org/10.1177/20539517221082027>.
- 1191 C. Knight, Équilibre réfléchissant. (2023) ; <https://plato.stanford.edu/entries/reflective-equilibrium/>.
- 1192 P. Kalluri, Ne demandez pas si l'intelligence artificielle est bonne ou juste, demandez-vous comment elle déplace le pouvoir. *Nature* 583, 169 (2020) ; <https://doi.org/10.1038/d41586-020-02003-2>.
- 1193 R. Dobbe, T. Krendl Gilbert, Y. Mintz, Choix difficiles en intelligence artificielle. *Intelligence artificielle* 300, 103555 (2021) ; <https://doi.org/10.1016/j.artint.2021.103555>.
- 1194* S. Fort, B. Lakshminarayanan, Ensemble Everything Everywhere : agrégation multi-échelle pour les confrontations Robustesse, arXiv [cs.CV] (2024) ; <http://arxiv.org/abs/2408.05446>.
- 1195 A. Zou, L. Phan, J. Wang, D. Duenas, M. Lin, M. Andriushchenko, R. Wang, Z. Kolter, M. Fredrikson, D. Hendrycks, Améliorer l'alignement et la robustesse avec des disjoncteurs, arXiv [cs.LG] (2024) ; <http://arxiv.org/abs/2406.04313>.
- 1196 M. Williams, M. Carroll, A. Narang, C. Weisser, B. Murphy, A. Dragan, Sur la manipulation et la tromperie ciblées Lors de l'optimisation des LLM pour les commentaires des utilisateurs, arXiv [cs.LG] (2024) ; <http://arxiv.org/abs/2411.02306>.
- 1197 S. Arnesen, D. Rein, J. Michael, La formation de modèles linguistiques pour gagner des débats avec l'auto-jeu améliore la précision des juges, arXiv [cs.CL] (2024) ; <http://arxiv.org/abs/2409.16636>.
- 1198 Z. Kenton, NY Siegel, J. Kramar, J. Brown-Cohen, S. Albanie, J. Bulian, R. Agarwal, D. Lindner, Y. Tang, N. Goodman, R. Shah, « Sur la surveillance évolutive des LLM faibles jugeant les LLM forts » dans 38e Conférence annuelle sur les systèmes de traitement de l'information neuronale (NeurIPS 2024) (2024) ; <https://openreview.net/forum?id=O1fp9nVraj>.
- 1199 A. Zou, L. Phan, S. Chen, J. Campbell, P. Guo, R. Ren, A. Pan, X. Yin, M. Mazeika, A.-K. Dombrowski, S. Goel, N. Li, MJ Byun, Z. Wang, A. Mallen, S. Basart, S. Koyejo, ... D. Hendrycks, Ingénierie des représentations : une approche descendante de la transparence de l'IA, arXiv [cs.LG] (2023) ; <http://arxiv.org/abs/2310.01405>.
- 1200 S. Casper, L. Schulze, O. Patel, D. Hadfield-Menell, Défense contre les modes de défaillance imprévus avec des Formation contradictoire, arXiv [cs.CR] (2024) ; <http://arxiv.org/abs/2403.05030>.
- 1201 TR Shaham, S. Schwettmann, F. Wang, A. Rajaram, E. Hernandez, J. Andreas, A. Torralba, A Multimodal Agent d'interprétabilité automatisé (2024) ; <https://openreview.net/forum?id=mDw42ZanmE>.
- 1202* Z. Kenton, T. Everitt, L. Weidinger, I. Gabriel, V. Mikulik, G. Irving, « Alignement des agents linguistiques » (Google DeepMind, 2021) ; <http://arxiv.org/abs/2103.14659>.
- 1203* C. Burns, P. Izmailov, JH Kirchner, B. Baker, L. Gao, L. Aschenbrenner, Y. Chen, A. Ecoffet, M. Joglekar, J. Leike, I. Sutskever, J. Wu, Généralisation faible-forte : susciter des capacités fortes avec une supervision faible, arXiv [cs.CL] (2023) ; <http://arxiv.org/abs/2312.09390>.
- 1204* J. Michael, S. Mahdi, D. Rein, J. Petty, J. Dirani, V. Padmakumar, SR Bowman, Le débat aide à superviser des Experts, arXiv [cs.AI] (2023) ; <http://arxiv.org/abs/2311.08702>.
- 1205 Y. Bengio, MK Cohen, N. Malkin, M. MacDermott, D. Fornasiere, P. Greiner, Y. Kaddar, Un oracle bayésien peut-il empêcher un agent de nuire ?, arXiv [cs.AI] (2024) ; <http://arxiv.org/abs/2408.05284>.
- 1206 M. Wu, AF Aji, Le style avant le fond : biais d'évaluation pour les grands modèles linguistiques, arXiv [cs.CL] (2023) ; <http://arxiv.org/abs/2307.03025>.
- 1207* N. Lambert, R. Calandra, Le plafond d'alignement : inadéquation objective dans l'apprentissage par renforcement de l'humain Commentaires, arXiv [cs.LG] (2023) ; <http://arxiv.org/abs/2311.00168>.
- 1208 H. Bansal, J. Dang, A. Grover, « Peering Through Preferences : Décrypter l'acquisition de rétroaction pour aligner les grands modèles linguistiques » dans la 12e Conférence internationale sur les représentations d'apprentissage (ICLR 2024) (Vienne, Autriche, 2023) ; <https://openreview.net/forum?id=dKI6IMwbCy>.
- 1209* J. Uesato, N. Kushman, R. Kumar, F. Song, N. Siegel, L. Wang, A. Creswell, G. Irving, I. Higgins, « Résoudre des problèmes mathématiques avec un retour d'information basé sur le processus et les résultats » (Google Deepmind, 2022) ; <https://doi.org/10.48550/arXiv.2211.14275>.
- 1210 H. Lightman, V. Kosaraju, Y. Burda, H. Edwards, B. Baker, T. Lee, J. Leike, J. Schulman, I. Sutskever, K. Cobbe, « Vérifions étape par étape » dans la 12e Conférence internationale sur les représentations d'apprentissage (ICLR 2024) (Vienne, Autriche, 2023) ; <https://openreview.net/forum?id=v8L0pN6EOi>.
- 1211 Z. Wu, Y. Hu, W. Shi, N. Dziri, A. Suhr, P. Ammanabrolu, NA Smith, M. Ostendorf, H. Hajishirzi, « Un retour d'information humain précis offre de meilleures récompenses pour l'entraînement au modèle linguistique » dans 37e Conférence sur les systèmes de traitement de l'information neuronale (NeurIPS 2023) (La Nouvelle-Orléans, LA, États-Unis, 2023) ; <https://openreview.net/forum?id=CSbGXyCswu>.
- 1212 Z. Li, Le côté obscur de ChatGPT : défis juridiques et éthiques posés par les perroquets stochastiques et les hallucinations, arXiv

- [cs.CY] (2023) ; <http://arxiv.org/abs/2304.14347>.
- 1213* A. Askeff, Y. Bai, A. Chen, D. Drain, D. Ganguli, T. Henighan, A. Jones, N. Joseph, B. Mann, N. DasSarma, N. Elhage, Z. Hatfield-Dodds, D. Hernandez, J. Kernion, K. Ndousse, C. Olsson, D. Amodei, ... J. Kaplan, Un assistant de langage général comme laboratoire d'alignement, arXiv [cs.CL] (2021) ; <http://arxiv.org/abs/2112.00861>.
- 1214 K. Shuster, S. Poff, M. Chen, D. Kiela, J. Weston, « L'augmentation de la récupération réduit les hallucinations dans les conversations » dans les conclusions de l'Association pour la linguistique computationnelle : EMNLP 2021, M.-F. Moens, X. Huang, L. Specia, SW-T. Yih, éd. (Association de linguistique computationnelle, Punta Cana, République dominicaine, 2021), pp. 3784–3803 ; <https://doi.org/10.18653/v1/2021.findings-emnlp.320>.
- 1215 L. Kuhn, Y. Gal, S. Farquhar, « Incertitude sémantique : invariances linguistiques pour l'estimation de l'incertitude dans la nature « Génération de langage » dans la 11e Conférence internationale sur les représentations de l'apprentissage (ICLR 2023) (Kigali, Rwanda, 2023) ; <https://openreview.net/forum?id=VD-AYtP0dve>.
- 1216 S. Min, K. Krishna, X. Lyu, M. Lewis, W.-T. Yih, P. Koh, M. Iyyer, L. Zettlemoyer, H. Hajishirzi, « FActScore : Fine-Évaluation atomique granulaire de la précision factuelle dans la génération de texte long » dans les actes de la conférence 2023 sur les méthodes empiriques en traitement du langage naturel (Association for Computational Linguistics, Stroudsburg, PA, États-Unis, 2023), pp. 12076–12100 ; <https://doi.org/10.18653/v1/2023.emnlp-main.741>.
- 1217 L. Chen, A. Perez-Lebel, FM Suchanek, G. Varoquaux, Reconfiguring LLMs from the Grouping Loss Perspective, arXiv [cs.CL] (2024) ; <http://arxiv.org/abs/2402.04957>.
- 1218 D. Hendrycks, S. Basart, N. Mu, S. Kadavath, F. Wang, E. Dorundo, R. Desai, T. Zhu, S. Parajuli, M. Guo, D. Song, J. Steinhardt, J. Gilmer, « Les multiples facettes de la robustesse : une analyse critique de la généralisation hors distribution » dans Conférence internationale IEEE/CVF 2021 sur la vision par ordinateur (ICCV) (2021), pp. 8320–8329 ; <https://doi.org/10.1109/ICCV48922.2021.00823>.
- 1219* S. Kadavath, T. Conerly, A. Askeff, T. Henighan, D. Drain, E. Pérez, N. Schiefer, Z. Hatfield-Dodds, N. DasSarma, E. Tran-Johnson, S. Johnston, S. El-Showk, A. Jones, N. Elhage, T. Hume, A. Chen, Y. Bai, ... J. Kaplan, Les modèles linguistiques (pour la plupart) savent ce qu'ils savent, arXiv [cs.CL] (2022) ; <http://arxiv.org/abs/2207.05221>.
- 1220* YA Yadkori, I. Kuzborskij, A. György, C. Szepesvári, Croire ou ne pas croire votre LLM, arXiv [cs.LG] (2024) ; <http://arxiv.org/abs/2406.02543>.
- 1221 S. Marks, C. Rager, EJ Michaud, Y. Belinkov, D. Bau, A. Mueller, Sparse Feature Circuits : découverte et édition de graphes causaux interprétables dans les modèles de langage, arXiv [cs.LG] (2024) ; <http://arxiv.org/abs/2403.19647>.
- 1222* T. Lieberum, M. Rahtz, J. Kramár, N. Nanda, G. Irving, R. Shah, V. Mikulik, « L'interprétabilité de l'analyse des circuits est-elle évolutive ? Preuves des capacités à choix multiples chez Chinchilla » (Google Deepmind, 2023) ; <https://doi.org/10.48550/arXiv.2307.09458>.
- 1223 E. Mitchell, C. Lin, A. Bosselut, CD Manning, C. Finn, « Édition de modèles basée sur la mémoire à grande échelle » dans Actes de la 39e Conférence internationale sur l'apprentissage automatique (PMLR, 2022), pp. 15817–15831 ; <https://proceedings.mlr.press/v162/mitchell22a.html>.
- 1224 K. Meng, AS Sharma, AJ Andonian, Y. Belinkov, D. Bau, « Édition de masse de la mémoire dans un transformateur » dans 11e Conférence internationale sur les représentations de l'apprentissage (ICLR 2023) (Kigali, Rwanda, 2022) ; <https://openreview.net/forum?id=MkbcAHlYgyS>.
- 1225 Y. Gandelsman, AA Efros, J. Steinhardt, « Interprétation de la représentation d'image de CLIP via une approche textuelle « Décomposition » dans la 12e Conférence internationale sur les représentations d'apprentissage (ICLR 2024) (Vienne, Autriche, 2023) ; <https://openreview.net/forum?id=5Ca9sSzuDp>.
- 1226 C. Tan, G. Zhang, J. Fu, « Édition massive pour les grands modèles de langage via le méta-apprentissage » dans The 12th International Conference on the Representations of Learning (ICLR 2024) (Vienne, Autriche, 2023) ; <https://openreview.net/forum?id=L6L1CJQ2PE>.
- 1227 S. Wang, Y. Zhu, H. Liu, Z. Zheng, C. Chen, J. Li, Édition des connaissances pour les grands modèles linguistiques : une étude, arXiv [cs.CL] (2023) ; <http://arxiv.org/abs/2310.16218>.
- 1228 A. Ghorbani, JY Zou, « Neuron Shapley : découverte des neurones responsables » dans Advances in Neural Information Processing Systems (NeurIPS 2020) (Curran Associates, Inc., 2020) vol. 33, pp. 5922–5932 ; <https://proceedings.neurips.cc/paper/2020/hash/41c542dfe6e4fc3deb251d64cf6ed2e4-Abstract.html>.
- 1229 X. Wu, J. Li, M. Xu, W. Dong, S. Wu, C. Bian, D. Xiong, « DEPN : détection et modification des neurones de confidentialité dans les modèles de langage pré-entraînés » dans Actes de la Conférence 2023 sur les méthodes empiriques en traitement du langage naturel (EMNLP 2023), H. Bouamor, J. Pino, K. Bali, éd. (Association for Computational Linguistics, Gateway, Singapour, 2023), pp. 2875–2886 ; <https://doi.org/10.18653/v1/2023.emnlp-main.174>.
- 1230 K. Li, O. Patel, F. Viégas, H. Pfister, M. Wattenberg, « Intervention au moment de l'inférence : susciter des réponses véridiques à partir d'un modèle linguistique » dans la 37e Conférence sur les systèmes de traitement de l'information neuronale (NeurIPS 2023) (La Nouvelle-Orléans,

- LA, États-Unis, 2023) ; <https://openreview.net/forum?id=aLLuYpn83y>.
- 1231 N. Belrose, D. Schneider-Joseph, S. Ravfogel, R. Cotterell, E. Raff, S. Biderman, « LEACE : effacement parfait du concept linéaire sous forme fermée » dans 37e Conférence sur les systèmes de traitement de l'information neuronale (NeurIPS 2023) (La Nouvelle-Orléans, LA, États-Unis, 2023) ; <https://openreview.net/forum?id=awlpKpwTWF¬elD=Ju4XcafMir>.
- 1232 AM Turner, L. Thiergart, D. Udell, G. Leech, U. Mini, M. MacDiarmid, Ajout d'activation : modèles de langage de pilotage sans optimisation, arXiv [cs.CL] (2023) ; <http://arxiv.org/abs/2308.10248>.
- 1233 E. Hernandez, BZ Li, J. Andreas, Inspection et édition des représentations de connaissances dans les modèles de langage, arXiv [cs.CL] (2023) ; <http://arxiv.org/abs/2304.00740>.
- 1234 D. Brown, C. Godfrey, C. Nizinski, J. Tu, H. Kvinge, « Robustesse des réseaux neuronaux édités » dans Atelier ICLR 2023 sur la compréhension mathématique et empirique des modèles de fondation (ME-FoMo 2023) (Kigali, Rwanda, 2023) ; <https://openreview.net/forum?id=JAjH6VANZ4>.
- 1235* C. Anil, E. Durmus, M. Sharma, J. Benton, S. Kundu, J. Batson, N. Rimsky, M. Tong, J. Mu, D. Ford, F. Mosconi, R. Agrawal, R. Schaeffer, N. Bashkanky, S. Svenningsen, M. Lambert, A. Radhakrishnan, ... D. Duvenaud, « Jailbreaking à plusieurs coups » (Anthropic, 2024) ; https://www-cdn.anthropic.com/af5633c94ed2beb282f6a53c595eb437e8e7b630/Many_Shot_Jailbreaking__2024_04_02_0936.pdf.
- 1236 Y. Deng, W. Zhang, SJ Pan, L. Bing, « Défis du jailbreak multilingue dans les grands modèles linguistiques » dans 12e Conférence internationale sur les représentations d'apprentissage (2024) ; <https://openreview.net/forum?id=vESNKdEMGp>.
- 1237 Y. Yuan, W. Jiao, W. Wang, J.-T. Huang, P. He, S. Shi, Z. Tu, « GPT-4 est trop intelligent pour être sûr : conversation furtive avec les LLM via Cipher » dans 12e Conférence internationale sur les représentations d'apprentissage (2024) ; <https://openreview.net/forum?id=MbFAK4s61A>.
- 1238 P. Ding, J. Kuang, D. Ma, X. Cao, Y. Xian, J. Chen, S. Huang, « Un loup déguisé en mouton : une approche imbriquée généralisée » « Les invites de jailbreak peuvent facilement tromper les grands modèles de langage » dans le chapitre nord-américain de l'Association for Computational Linguistics (2023) ; <https://api.semanticscholar.org/CorpusID:265664913>.
- 1239 Z. Wei, Y. Wang, A. Li, Y. Mo, Y. Wang, Jailbreak et modèles de langage alignés sur la garde avec seulement quelques éléments en contexte Démonstrations, arXiv [cs.LG] (2023) ; <http://arxiv.org/abs/2310.06387>.
- 1240* M. Russinovich, A. Salem, R. Eldan, Super, écrivez maintenant un article à ce sujet : Le LLM multi-tours Crescendo Attaque de jailbreak, arXiv [cs.CR] (2024) ; <http://arxiv.org/abs/2404.01833>.
- 1241 A. Madry, A. Makelov, L. Schmidt, D. Tsipras, A. Vladu, « Vers des modèles d'apprentissage profond résistants aux conflits « Attacks » dans la 6e Conférence internationale sur les représentations de l'apprentissage (ICLR 2018) (Vancouver, BC, Canada, 2018) ; <https://openreview.net/forum?id=rJzIBfZAb>.
- 1242 S. Friedler, R. Singh, B. Bili-Hamelin, J. Metcalf, BJ Chen, « Le red-teaming de l'IA n'est pas une solution unique à l'IA Harms : Recommandations pour l'utilisation du Red-Teaming pour la responsabilisation de l'IA (Data & Society, 2023) ; <https://datasociety.net/library/ai-red-teaming-is-not-a-one-stop-solution-to-ai-harms-recommendations-for-using-red-teaming-for-ai-accountability/>.
- 1243 N. Jain, A. Schwarzschild, Y. Wen, G. Somepalli, J. Kirchenbauer, P.-Y. Chiang, M. Goldblum, A. Saha, J. Geiping, T. Goldstein, Défenses de base pour les attaques adverses contre les modèles de langage alignés, arXiv [cs.LG] (2023) ; <http://arxiv.org/abs/2309.00614>.
- 1244 S. Lee, M. Kim, L. Cherif, D. Dobre, J. Lee, SJ Hwang, K. Kawaguchi, G. Gidel, Y. Bengio, N. Malkin, M. Jain, Apprentissage de diverses attaques sur de grands modèles de langage pour un Red-Teaming robuste et un réglage de sécurité, arXiv [cs.CL] (2024) ; <http://arxiv.org/abs/2405.18540>.
- 1245 A. Peng, J. Michael, H. Sleight, E. Perez, M. Sharma, Réponse rapide : atténuer les jailbreaks LLM avec quelques exemples, arXiv [cs.CL] (2024) ; <http://arxiv.org/abs/2411.07494>.
- 1246 Z. Liu, G. Dou, Z. Tan, Y. Tian, M. Jiang, Vers des modèles de langage plus sûrs et de grande taille grâce au désapprentissage automatique, arXiv [cs.CL] (2024) ; <http://arxiv.org/abs/2402.10058>.
- 1247 A. Lynch, P. Guo, A. Ewart, S. Casper, D. Hadfield-Menell, Huit méthodes pour évaluer le désapprentissage robuste dans les LLM, arXiv [cs.CL] (2024) ; <http://arxiv.org/abs/2402.16835>.
- 1248 D. Gamage, J. Chen, K. Sasahara, « L'émergence des deepfakes et ses implications sociétales : une étude systématique Critique » dans la Conférence pour la vérité et la confiance en ligne 2021 (2021), pp. 28–39 ; https://www.researchgate.net/publication/355583941_The_Emergence_of_Deepfakes_and_its_Societal_Implications_Une_revue_systématique.
- 1249 A. Kaushal, A. Mina, A. Meena, TH Babu, « L'impact sociétal des deepfakes : progrès dans la détection et « Mitigation » dans la 14e Conférence internationale sur les technologies informatiques, de communication et de réseau (ICCCNT) (2023), pp. 1–7 ; <https://doi.org/10.1109/ICCCNT56998.2023.10307353>.

- 1250 F. Romero Moreno, IA générative et Deepfakes : une approche fondée sur les droits de l'homme pour lutter contre les contenus préjudiciables. *Revue internationale de droit Informatique & Technologie* 38, 297–326 (2024) ; <https://doi.org/10.1080/13600869.2024.2324540>.
- 1251 R. Tang, Y.-N. Chuang, X. Hu, La science de la détection de texte généré par LLM. *Communications de l'ACM* 67, 50–59 (04/2024) ; <https://doi.org/10.1145/3624725>.
- 1252 K. Krishna, Y. Song, M. Karpinska, JF Wieting, M. Iyyer, « La paraphrase échappe aux détecteurs de texte généré par l'IA, mais la récupération est une défense efficace » dans 37e Conférence sur les systèmes de traitement de l'information neuronale (NeurIPS 2023) (2023) ; <https://openreview.net/pdf?id=WbFhFvjjKj>.
- 1253 L. Lin, N. Gupta, Y. Zhang, H. Ren, C.-H. Liu, F. Ding, X. Wang, X. Li, L. Verdoliva, S. Hu, Détection du multimédia généré par de grands modèles d'IA : une enquête, *arXiv [cs.MM]* (2024) ; <http://arxiv.org/abs/2402.00045>.
- 1254 R. Corvi, D. Cozzolino, G. Zingarini, G. Poggi, K. Nagano, L. Verdoliva, « Sur la détection d'images synthétiques « Généré par des modèles de diffusion » dans ICASSP 2023 - Conférence internationale IEEE 2023 sur l'acoustique, la parole et le traitement du signal (ICASSP) (2023), pp. 1–5 ; <https://doi.org/10.1109/ICASSP49357.2023.10095167>.
- 1255 U. Ojha, Y. Li, YJ Lee, « Vers des détecteurs universels de fausses images qui se généralisent à travers les modèles génératifs » dans Conférence IEEE/CVF 2023 sur la vision par ordinateur et la reconnaissance de formes (CVPR) (IEEE Computer Society, 2023), pp. 24480–24489 ; <https://doi.org/10.1109/CVPR52729.2023.02345>.
- 1256 HB Wee, JD Reimer, Les universitaires non anglophones confrontés aux inégalités via des essais générés par l'IA et des contre-mesures Outils. *Bioscience* 73, 476–478 (2023) ; <https://doi.org/10.1093/biosci/biad034>.
- 1257 Y. Zhao, T. Pang, C. Du, X. Yang, N.-M. Cheung, M. Lin, Une recette pour les modèles de diffusion de filigrane, *arXiv [cs.CV]* (2023) ; <http://arxiv.org/abs/2303.10137>.
- 1258 M. Christ, S. Gunn, O. Zamir, « Filigranes indétectables pour les modèles linguistiques » dans les actes de la 37e Conférence sur la théorie de l'apprentissage, S. Agrawal, A. Roth, éd. (PMLR, 2024) vol. 247 des Actes de la recherche sur l'apprentissage automatique, pp. 1125–1139 ; <https://proceedings.mlr.press/v247/christ24a.html>.
- 1259 J. Kirchenbauer, J. Geiping, Y. Wen, J. Katz, I. Miers, T. Goldstein, « Un filigrane pour les grands modèles de langage » dans Actes de la 40e Conférence internationale sur l'apprentissage automatique (PMLR, 2023), pp. 17061–17084 ; <https://proceedings.mlr.press/v202/kirchenbauer23a.html>.
- 1260 Y. Liu, Y. Bu, « Filigrane de texte adaptatif pour les grands modèles linguistiques » dans la quarante-et-unième conférence internationale sur l'apprentissage automatique (2024) ; <https://openreview.net/forum?id=7emOSb5UfX>.
- 1261 A. Liu, L. Pan, Y. Lu, J. Li, X. Hu, X. Zhang, L. Wen, I. King, H. Xiong, PS Yu, Une étude du filigranage de texte à l'époque des grands modèles de langage, *arXiv [cs.CL]* (2023) ; <http://arxiv.org/abs/2312.07913>.
- 1262 H. Zhang, BL Edelman, D. Francati, D. Venturi, G. Ateniese, B. Barak, Filigranes dans le sable : impossibilité de Filigranage fort pour les modèles génératifs, *arXiv [cs.LG]* (2023) ; <http://arxiv.org/abs/2311.04378>.
- 1263 A. Knott, D. Pedreschi, R. Chatila, T. Chakraborti, S. Leavy, R. Baeza-Yates, D. Eysers, A. Trotman, PD Teal, P. Biecek, S. Russell, Y. Bengio, Les modèles d'IA génératifs devraient inclure des mécanismes de détection comme condition de publication publique. *Éthique et technologies de l'information* 25, 55 (2023) ; <https://doi.org/10.1007/s10676-023-09728-4>.
- 1264 C2PA, Aperçu (2022) ; <https://c2pa.org/>.
- 1265 AI for Good, collaboration sur les normes d'authenticité de l'IA et du multimédia (2024) ; <https://aiforgood.itu.int/multimedia-authenticity/>.
- 1266 A. Al-Dhaqm, RA Ikuesan, VR Kebande, SA Razak, G. Grispos, K.-KR Choo, BAS Al-Rimy, AA Alsewari, Sous-domaines de la criminalistique numérique : l'état de l'art et les orientations futures. *IEEE Access* 9, 152476–152502 (2021) ; <https://doi.org/10.1109/ACCESS.2021.3124262>.
- 1267 F. Casino, TK Dasaklis, GP Spathoulas, M. Anagnostopoulos, A. Ghosal, I. Borocz, A. Solanas, M. Conti, C. Patsakis, Tendances de recherche, défis et sujets émergents en criminalistique numérique : une revue des revues. *IEEE Access* 10, 25464–25493 (2022) ; <https://doi.org/10.1109/ACCESS.2022.3154059>.
- 1268 HR Hasan, K. Salah, Combattre les vidéos deepfake à l'aide de la blockchain et des contrats intelligents. *IEEE Access : pratique Innovations, solutions ouvertes* 7, 41596–41606 (2019) ; <https://doi.org/10.1109/access.2019.2905689>.
- 1269 CC Ki Chan, V. Kumar, S. Delaney, M. Gochoo, « Combattre les Deepfakes : Multi-LSTM et Blockchain comme preuve d'authenticité pour les médias numériques » dans Conférence internationale IEEE/ITU 2020 sur l'intelligence artificielle pour le bien (AI4G) (IEEE, 2020) ; <https://doi.org/10.1109/ai4g50087.2020.9311067>.
- 1270 P. Fraga-Lamas, TM Fernández-Caramés, Fake News, Disinformation, and Deepfakes: Tirer parti des technologies de registre distribué et de la blockchain pour lutter contre la tromperie numérique et la réalité contrefaite, *arXiv [cs.CY]* (2019) ; <http://dx.doi.org/10.1109/MITP.2020.2977589>.
- 1271 S. Mohammad Niyaz Khan, J. Mohd Ghazali, LQ Zakaria, SN Ahmad, KA Elias, diverses classifications d'images Utilisation de certaines métadonnées EXIF (Exchangeable Image File Format) des images. *Revue malaisienne d'information et*

- Technologies de la communication (MyJICT), 1–12 (2018) ; <https://doi.org/10.53840/myjict3-1-33>.
- 1272 A. Chan, C. Ezell, M. Kaufmann, K. Wei, L. Hammond, H. Bradley, E. Bluemke, N. Rajkumar, D. Krueger, N. Kolt, L. Heim, M. Anderljung, « Visibilité sur les agents d'IA » dans la conférence ACM 2024 sur l'équité, la responsabilité et la transparence (ACM, New York, NY, États-Unis, 2024) ; <https://doi.org/10.1145/3630106.3658948>.
- 1273 A. Chan, N. Kolt, P. Wills, U. Anwar, C. S. de Witt, N. Rajkumar, L. Hammond, D. Krueger, L. Heim, M. Anderljung, IDs pour les systèmes d'IA, arXiv [cs.AI] (2024) ; <http://arxiv.org/abs/2406.12137>.
- 1274 B. Pan, N. Stakhanova, S. Ray, Provenance des données en matière de sécurité et de confidentialité. ACM Computing Surveys 55, 1–35 (2023) ; <https://doi.org/10.1145/3593294>.
- 1275 E. Laird, M. Dwyer, « Hors tâche : menaces EdTech pour la vie privée et l'équité des étudiants à l'ère de l'IA » (Center for Démocratie et technologie, 2023) ; <https://cdt.org/insights/report-off-task-edtech-threats-to-student-privacy-and-equity-in-the-age-of-ai/>.
- 1276 SS El Mokadem, L'effet de l'éducation aux médias sur la désinformation et la détection de fausses vidéos. Médias arabes et Société (2023) ; <https://www.arabmediasociety.com/the-effect-of-media-literacy-on-misinformation-and-deep-fake-video-detection/>.
- 1277 Y. Hwang, JY Ryu, S.-H. Jeong, Effets de la désinformation utilisant le deepfake : l'effet protecteur des médias Éducation à la littératie. Cyberpsychologie, comportement et réseaux sociaux 24, 188–193 (2021) ; <https://doi.org/10.1089/cyber.2020.0174>.
- 1278 SY Shin, J. Lee, L'effet des vidéos deepfakes sur la crédibilité des informations et l'influence corrective des connaissances basées sur les coûts concernant les deepfakes. Digital Journalism 10, 412–432 (2022) ; <https://doi.org/10.1080/21670811.2022.2026797>.
- 1279 S. Qian, C. Shen, J. Zhang, Combattre les contrefaçons bon marché : utiliser une intervention d'éducation aux médias numériques pour motiver la recherche inversée de fausses informations visuelles hors contexte. Journal of Computer-Mediated Communication : JCMC 28 (2022) ; <https://doi.org/10.1093/jcmc/zmac024>.
- 1280 T. Ali, P. Kostakos, HuntGPT : Intégration de la détection d'anomalies basée sur l'apprentissage automatique et de l'IA explicable avec de grands modèles de langage (LLM), arXiv [cs.CR] (2023) ; <http://arxiv.org/abs/2309.16021>.
- 1281 G. Pang, C. Shen, L. Cao, A. Van Den Hengel, Apprentissage profond pour la détection des anomalies : une revue. ACM Computing Enquêtes 54, 38 : 1–38 : 38 (2021) ; <https://doi.org/10.1145/3439950>.
- 1282 J. Geng, F. Cai, Y. Wang, H. Koepl, P. Nakov, I. Gurevych, A Survey of Confidence Estimation and Calibration in Grands modèles de langage, arXiv [cs.CL] (2023) ; <http://arxiv.org/abs/2311.08298>.
- 1283 A. Aldahdooh, W. Hamidouche, SA Fezza, O. Déforges, Détection d'exemples contradictoires pour les modèles DNN : revue et comparaison expérimentale. Artificial Intelligence Review 55, 4403–4462 (2022) ; <https://doi.org/10.1007/s10462-021-10125-w>.
- 1284 J. Hayase, W. Kong, R. Somani, S. Oh, « SPECTRE : se défendre contre les attaques par porte dérobée à l'aide de statistiques robustes » dans les actes de la 38e Conférence internationale sur l'apprentissage automatique, M. Meila, T. Zhang, éd. (PMLR, 2021) vol. 139 des actes de recherche sur l'apprentissage automatique, pp. 4129–4139 ; <https://proceedings.mlr.press/v139/hayase21a.html>.
- 1285 AT Mallen, N. Belrose, « Élicitation de connaissances latentes à partir de modèles linguistiques originaux » dans l'atelier ICLR 2024 sur Compréhension mathématique et empirique des modèles de fondation (2024) ; <https://openreview.net/forum?id=Z1531QeqAQ>.
- 1286* M. MacDiarmid, T. Maxwell, N. Schiefer, J. Mu, J. Kaplan, D. Duvenaud, S. Bowman, A. Tamkin, E. Perez, M. Sharma, C. Denison, E. Hubinger, De simples sondes peuvent attraper des agents dormants (2024) ; <https://www.anthropic.com/news/probes-catch-sleeper-agents>.
- 1287 S. Han, K. Rao, A. Ettinger, L. Jiang, BY Lin, N. Lambert, Y. Choi, N. Dziri, « WildGuard : outils de modération ouverts et uniques pour les risques de sécurité, les jailbreaks et les refus de LLM » dans 38e conférence sur les ensembles de données et les repères des systèmes de traitement de l'information neuronale (2024) ; <https://openreview.net/forum?id=lch4tv4202>.
- 1288 R. Greenblatt, B. Shlegeris, K. Sachan, F. Roger, AI Control : améliorer la sécurité malgré la subversion intentionnelle, arXiv [cs.LG] (2023) ; <http://arxiv.org/abs/2312.06942>.
- 1289 M. Phute, A. Helbling, MD Hull, S. Peng, S. Szyller, C. Cornelius, DH Chau, « LLM Self Defense : en s'examinant, les LLM savent qu'ils sont trompés » dans The Second Tiny Papers Track à l'ICLR 2024 (Vienne, Autriche, 2024) ; <https://openreview.net/forum?id=YogqclA19o>.
- 1290* H. Inan, K. Upasani, J. Chi, R. Rungta, K. Iyer, Y. Mao, M. Tontchev, Q. Hu, B. Fuller, D. Testuggine, M. Khabsa, Llama Guard : LLM - Sauvegarde d'entrée-sortie basée sur les conversations homme-IA, arXiv [cs.CL] (2023) ; <http://arxiv.org/abs/2312.06674>.
- 1291 T. Kim, S. Kotha, A. Raghunathan, Jailbreaking des défenses avec le problème violet, arXiv [cs.CR] (2024) ;

- <http://arxiv.org/abs/2403.14725>.
- 1292 SO Hansson, M.-Å. Belin, B. Lundgren, Véhicules autonomes : aperçu éthique. *Philosophie et technologie* 34, 1383-1408 (2021) ; <https://doi.org/10.1007/s13347-021-00464-5>.
- 1293 NR Jennings, L. Moreau, D. Nicholson, S. Ramchurn, S. Roberts, T. Rodden, A. Rogers, Collectifs humains-agents. *Communications de l'ACM* 57, 80–88 (2014) ; <https://doi.org/10.1145/2629559>.
- 1294* A. Dafoe, E. Hughes, Y. Bachrach, T. Collins, KR McKee, JZ Leibo, K. Larson, T. Graepel, Problèmes ouverts dans IA coopérative, *arXiv [cs.AI]* (2020) ; <http://arxiv.org/abs/2012.08630>.
- 1295 A. Dafoe, Y. Bachrach, G. Hadfield, E. Horvitz, K. Larson, T. Graepel, IA coopérative : les machines doivent apprendre à trouver un terrain d'entente. *Nature* 593, 33–36 (2021) ; <https://doi.org/10.1038/d41586-021-01170-0>.
- 1296 D. Hadfield-Menell, A. Dragan, P. Abbeel, S. Russell, « Apprentissage coopératif par renforcement inverse » dans *Actes de la 30e Conférence internationale sur les systèmes de traitement de l'information neuronale (NIPS 2016)* (Curran Associates Inc., Red Hook, NY, États-Unis, 2016), pp. 3916–3924 ; https://papers.nips.cc/paper_files/paper/2016/hash/c3395dd46c34fa7fd8d729d8cf88b7a8-Abstract.html.
- 1297 I. Seeber, E. Bittner, RO Briggs, T. de Vreede, G.-J. de Vreede, A. Elkins, R. Maier, AB Merz, S. Oeste-Reiß, N. Randrup, G. Schwabe, M. Söllner, Les machines comme coéquipiers : un programme de recherche sur l'IA dans la collaboration en équipe. *Information & Gestion* 57, 103174 (2020) ; <https://doi.org/10.1016/j.im.2019.103174>.
- 1298 R. Shah, P. Freire, N. Alex, R. Freedman, D. Krashennikov, L. Chan, MD Dennis, P. Abbeel, A. Dragan, S. Russell, Avantages de l'assistance par rapport à l'apprentissage par récompense (2020) ; <https://openreview.net/forum?id=DFIoGDZejlB>.
- 1299 SD Ramchurn, S. Stein, NR Jennings, Partenariats homme-IA dignes de confiance. *iScience* 24, 102891 (2021) ; <https://doi.org/10.1016/j.isci.2021.102891>.
- 1300 X. Wu, L. Xiao, Y. Sun, J. Zhang, T. Ma, L. He, Une étude sur l'implication humaine dans la boucle pour l'apprentissage automatique. *Future Generations Computer Systems : FGCS* 135, 364–381 (2022) ; <https://doi.org/10.1016/j.future.2022.05.014>.
- 1301 KL Mosier, LJ Skitka, Utilisation de l'automatisation et biais d'automatisation. *Actes de la conférence Human Factors and Société d'ergonomie... Réunion annuelle. Société des facteurs humains et de l'ergonomie. Réunion annuelle* 43, 344–348 (1999) ; <https://doi.org/10.1177/154193129904300346>.
- 1302 J. Babcock, J. Krámar, RV Yampolskiy, « Lignes directrices pour le confinement de l'intelligence artificielle » dans *Next-Generation Éthique : concevoir une société meilleure*, AE Abbas, éd. (Cambridge University Press, Cambridge, 2019), pp. 90–112 ; <https://doi.org/10.1017/9781108616188.008>.
- 1303 SG Patil, T. Zhang, V. Fang, NC, R. Huang, A. Hao, M. Casado, JE Gonzalez, RA Popa, I. Stoica, GoEX : perspectives et conceptions vers un environnement d'exécution pour les applications LLM autonomes, *arXiv [cs.CL]* (2024) ; <http://arxiv.org/abs/2404.06921>.
- 1304 J. Gryz, M. Rojszczak, Algorithmes de boîte noire et droits des individus : pas de solution facile à l'« explicabilité » Problème. *Internet Policy Review* 10 (2021) ; <https://policyreview.info/articles/analysis/black-box-algorithms-and-rights-individuals-no-easy-solution-explainability>.
- 1305 JA McDermid, Y. Jia, Z. Porter, I. Habli, Explicabilité de l'intelligence artificielle : les dimensions techniques et éthiques. *Transactions philosophiques. Série A, Sciences mathématiques, physiques et de l'ingénierie* 379, 20200363 (2021) ; <https://doi.org/10.1098/rsta.2020.0363>.
- 1306 T. Ploug, S. Holm, « Droit de contester les diagnostics d'IA définissant les exigences de transparence et d'explicabilité du point de vue d'un patient » dans *Artificial Intelligence in Medicine* (Springer Publishing Company, 2022), pp. 227–238 ; https://doi.org/10.1007/978-3-030-64573-1_267.
- 1307 SH Tanneru, D. Ley, C. Agarwal, H. Lakkaraju, Sur la difficulté du raisonnement fidèle par chaîne de pensée dans les grandes Modèles de langage, *arXiv [cs.CL]* (2024) ; <http://arxiv.org/abs/2406.10625>.
- 1308* J. Chua, E. Rees, H. Batra, SR Bowman, J. Michael, E. Perez, M. Turpin, L'entraînement à la cohérence augmentée par biais réduit le raisonnement biaisé dans la chaîne de pensée, *arXiv [cs.CL]* (2024) ; <http://arxiv.org/abs/2403.05518>.
- 1309* A. Radhakrishnan, K. Nguyen, A. Chen, C. Chen, C. Denison, D. Hernandez, E. Durmus, E. Hubinger, J. Kernion, K. Lukošiuūtė, N. Cheng, N. Joseph, N. Schiefer, O. Rausch, S. McCandlish, S. El Showk, T. Lanham, ... E. Perez, La décomposition des questions améliore la fidélité du raisonnement généré par un modèle, *arXiv [cs.CL]* (2023) ; <http://arxiv.org/abs/2307.11768>.
- 1310 J. Li, P. Cao, Y. Chen, K. Liu, J. Zhao, Vers une chaîne de pensée fidèle : les grands modèles linguistiques font le pont Raisonnements, *arXiv [cs.CL]* (2024) ; <http://arxiv.org/abs/2405.18915>.
- 1311 D. Paul, R. West, A. Bosselut, B. Faltings, Making Reasoning Matter : mesurer et améliorer la fidélité du raisonnement par chaîne de pensée, *arXiv [cs.CL]* (2024) ; <http://arxiv.org/abs/2402.13950>.
- 1312 A. Saranya, R. Subhashini, Une revue systématique des modèles et applications d'intelligence artificielle explicables : développements récents et tendances futures. *Decision Analytics Journal* 7, 100230 (2023) ;

- <https://doi.org/10.1016/j.dajour.2023.100230>.
- 1313 H. Zhao, H. Chen, F. Yang, N. Liu, H. Deng, H. Cai, S. Wang, D. Yin, M. Du, Explicabilité des grands modèles de langage : une étude. *ACM Transactions on Intelligent Systems and Technology* 15, 1–38 (2024) ; <https://doi.org/10.1145/3639372>.
- 1314 S. Casper, C. Ezell, C. Siegmann, N. Kolt, TL Curtis, B. Bucknall, A. Haupt, K. Wei, J. Scheurer, M. Hobbahn, L. Sharkey, S. Krishna, M. Von Hagen, S. Alberti, A. Chan, Q. Sun, M. Gerovitch, ... D. Hadfield-Menell, « L'accès à la boîte noire est insuffisant pour des audits d'IA rigoureux » dans la conférence 2024 de l'ACM sur l'équité, la responsabilité et la transparence (ACM, New York, NY, États-Unis, 2024), pp. 2254–2272 ; <https://doi.org/10.1145/3630106.3659037>.
- 1315 O. Aarne, T. Fist, C. Withers, « Puces sécurisées et gouvernables : utiliser des mécanismes sur puce pour gérer les risques de sécurité nationale liés à l'IA et à l'informatique avancée » (Center for a New American Security, 2024) ; <https://s3.us-east-1.amazonaws.com/files.cnas.org/documents/CNAS-Report-Tech-Secure-Chips-Jan-24-finalb.pdf> .
- 1316 G. Kulp, D. Gonzales, E. Smith, L. Heim, P. Puri, M. Vermeer, Z. Winkelman, « Mécanismes de gouvernance basés sur le matériel » (RAND Corporation, 2024) ; https://www.rand.org/pubs/working_papers/WRA3056-1.html.
- 1317 Z. Ghodsi, T. Gu, S. Garg, SafetyNets : exécution vérifiable de réseaux neuronaux profonds sur un cloud non fiable. *Progress dans les systèmes de traitement de l'information neuronale* 30 (2017) ; https://proceedings.neurips.cc/paper_files/paper/2017/file/6048ff4e8cb07aa60b6777b6f7384d52-Paper.pdf.
- 1318 H. Chen, C. Fu, BD Rouhani, J. Zhao, F. Koushanfar, « DeepAttest : un cadre d'attestation de bout en bout pour les réseaux neuronaux profonds » dans Actes du 46e Symposium international sur l'architecture informatique (Association for Computing Machinery, New York, NY, États-Unis, 2019)ISCA '19, pp. 487–498 ; <https://doi.org/10.1145/3307650.3322251>.
- 1319 H. Jia, M. Yaghini, CA Choquette-Choo, N. Dullerud, A. Thudi, V. Chandrasekaran, N. Papernot, « Proof-of-Apprentissage : définitions et pratiques » dans Symposium IEEE 2021 sur la sécurité et la confidentialité (SP) (IEEE, 2021), pp. 1039–1056 ; <https://doi.org/10.1109/SP40001.2021.00106>.
- 1320 S. Goldwasser, GN Rothblum, J. Shafer, A. Yehudayoff, « Preuves interactives pour la vérification de l'apprentissage automatique » dans la 12e conférence sur les innovations en informatique théorique (ITCS 2021), JR Lee, éd. (Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl, Allemagne, 2021) vol. 185 de Leibniz International Proceedings in Informatics (LIPIcs), pp. 41 :1–41 :19 ; <https://doi.org/10.4230/LIPIcs.ITCS.2021.41>.
- 1321* Apple, « Sécurité de la plateforme Apple » (Apple, 2024) ; https://help.apple.com/pdf/security/en_US/apple-platform-guide-de-securite.pdf.
- 1322* J. Zhu, H. Yin, P. Deng, A. Almeida, S. Zhou, Informatique confidentielle sur GPU nVIDIA H100 : une étude de performance Étude de référence, arXiv [cs.DC] (2024) ; <http://arxiv.org/abs/2409.03992>.
- 1323 R. Anderson, S. Fuloria, « Qui contrôle l'interrupteur d'arrêt ? » dans Première conférence internationale IEEE 2010 sur les communications de réseau intelligent (IEEE, 2010), pp. 96–101 ; <https://doi.org/10.1109/smartgrid.2010.5622026>.
- 1324 Organisation de coopération et de développement économiques, « Technologies émergentes améliorant la confidentialité » (OCDE, 2023) ; <https://doi.org/10.1787/bf121be4-fr>.
- 1325 N. Subramani, S. Luccioni, J. Dodge, M. Mitchell, « Détection d'informations personnelles dans les corpus de formation : une analyse » dans les actes du 3e atelier sur le traitement fiable du langage naturel (TrustNLP 2023), A. Ovalle, K.-W. Chang, N. Mehrabi, Y. Pruksachatkun, A. Galystan, J. Dhamala, A. Verma, T. Cao, A. Kumar, R. Gupta, Eds. (Association de linguistique computationnelle, Toronto, Canada, 2023), pp. 208–220 ; <https://doi.org/10.18653/v1/2023.trustnlp-1.18>.
- 1326 Y. Elazar, A. Bhagia, IH Magnusson, A. Ravichander, D. Schwenk, A. Suhr, EP Walsh, D. Groeneveld, L. Soldaini, S. Singh, H. Hajishirzi, NA Smith, J. Dodge, « Qu'y a-t-il dans mes Big Data ? » dans 12e Conférence internationale sur les représentations d'apprentissage (2024) ; <https://openreview.net/forum?id=RvfPnOkPV4>.
- 1327 A. Narayanan, V. Shmatikov, « Désanonymisation robuste de grands ensembles de données éparses » dans Symposium IEEE 2008 sur la sécurité et la confidentialité (sp 2008) (2008), pp. 111–125 ; <https://doi.org/10.1109/SP.2008.33>.
- 1328 H. Brown, K. Lee, F. Mireshghallah, R. Shokri, F. Tramèr, « Que signifie pour un modèle linguistique préserver Confidentialité ? » dans Actes de la conférence 2022 de l'ACM sur l'équité, la responsabilité et la transparence (FAccT '22) (Association for Computing Machinery, New York, NY, États-Unis, 2022), pp. 2280–2292 ; <https://doi.org/10.1145/3531146.3534642>.
- 1329* S. Wu, O. Irsoy, S. Lu, V. Dabravolski, M. Dredze, S. Gehrmann, P. Kambadur, D. Rosenberg, G. Mann, BloombergGPT : un grand modèle de langage pour la finance, arXiv [cs.LG] (2023) ; <http://arxiv.org/abs/2303.17564>.
- 1330 G. Penedo, Q. Malartic, D. Hesslow, R. Cojocar, H. Alobeidli, A. Cappelli, B. Pannier, E. Almazrouei, J. Launay, « L'ensemble de données Web raffiné pour Falcon LLM : surperformer les corpus organisés avec des données Web uniquement » dans le cadre de la 37e conférence sur les systèmes de traitement de l'information neuronale (NeurIPS 2023) Datasets and Benchmarks Track (La Nouvelle-Orléans, LA,

- États-Unis, 2023) ; <https://openreview.net/pdf?id=kM5eGcdCzq>.
- 1331 T. Gebru, J. Morgenstern, B. Vecchione, JW Vaughan, H. Wallach, HD Iii, K. Crawford, Fiches techniques pour ensembles de données. *Communications de l'ACM* 64, 86–92 (2021) ; <https://doi.org/10.1145/3458723>.
- 1332 A. Ghorbani, J. Zou, « Data Shapley : évaluation équitable des données pour l'apprentissage automatique » dans Actes de la 36e Conférence internationale sur l'apprentissage automatique (ICML 2019), K. Chaudhuri, R. Salakhutdinov, éd. (PMLR, La Nouvelle-Orléans, LA, États-Unis, 2019) vol. 97 des Actes de la recherche sur l'apprentissage automatique, pp. 2242–2251 ; <https://proceedings.mlr.press/v97/ghorbani19c.html>.
- 1333 T. Li, EF Villaronga, P. Kieseberg, Les humains oublient, les machines se souviennent : l'intelligence artificielle et le droit d'être Oublié. *Revue de droit et de sécurité informatique* 34, 304 (2018) ; https://scholarship.law.bu.edu/faculty_scholarship/817.
- 1334 Z. Zhang, M. Jia, H.-P. Lee, B. Yao, S. Das, A. Lerner, D. Wang, T. Li, « C'est un jeu équitable », ou pas ? Examen de la façon dont les utilisateurs naviguent entre les risques et les avantages de la divulgation lors de l'utilisation d'agents conversationnels basés sur LLM, *arXiv [cs.HC]* (2023) ; <http://dx.doi.org/10.1145/3613904.3642385>.
- 1335 Z. Zhang, C. Shen, B. Yao, D. Wang, T. Li, Utilisation secrète du grand modèle de langage (LLM), *arXiv [cs.HC]* (2024) ; <http://arxiv.org/abs/2409.19450>.
- 1336 C. Dwork, F. McSherry, K. Nissim, A. Smith, « Calibrage du bruit en fonction de la sensibilité dans l'analyse des données privées » dans *Theory of Cryptography*, S. Halevi, T. Rabin, éd. (Springer, Berlin, Heidelberg, 2006) vol. 3876 de *Lecture Notes in Computer Science* ; https://doi.org/10.1007/11681878_14.
- 1337 M. Abadi, A. Chu, I. Goodfellow, HB McMahan, I. Mironov, K. Talwar, L. Zhang, « Apprentissage profond avec calcul différentiel Confidentialité » dans les actes de la conférence ACM SIGSAC 2016 sur la sécurité informatique et des communications (CCS '16) (Association for Computing Machinery, New York, NY, États-Unis, 2016), pp. 308–318 ; <https://doi.org/10.1145/2976749.2978318>.
- 1338* S. De, L. Berrada, J. Hayes, SL Smith, B. Balle, « Déverrouillage d'une image différentiellement privée de haute précision « Classification par échelle » (Google Deepmind, 2022) ; <http://arxiv.org/abs/2204.13650>.
- 1339 X. Li, F. Tramer, P. Liang, T. Hashimoto, « Les grands modèles linguistiques peuvent être des apprenants différentiels forts et privés » dans Conférence internationale sur les représentations d'apprentissage 2022 (virtuelle, 2022) ; <https://openreview.net/forum?id=bVuP3ItATMz>.
- 1340 D. Yu, S. Naik, A. Backurs, S. Gopi, HA Inan, G. Kamath, J. Kulkarni, YT Lee, A. Manoel, L. Wutschitz, S. Yekhanin, H. Zhang, « Réglage fin différentiellement privé des modèles de langage » dans Conférence internationale sur les représentations d'apprentissage (2022) ; <https://openreview.net/forum?id=Q42f0dfjECO>.
- 1341* A. Kurakin, N. Ponomareva, U. Syed, L. MacDermed, A. Terzis, Exploiter les modèles à grand langage pour générer Texte synthétique privé, *arXiv [cs.LG]* (2023) ; <http://arxiv.org/abs/2306.01684>.
- 1342 R. Liu, J. Wei, F. Liu, C. Si, Y. Zhang, J. Rao, S. Zheng, D. Peng, D. Yang, D. Zhou, AM Dai, « Meilleures pratiques et leçons apprises sur les données synthétiques » dans Première conférence sur la modélisation du langage (2024) ; <https://openreview.net/forum?id=OJaWBhh61C>.
- 1343 A. Yale, S. Dash, R. Dutta, I. Guyon, A. Pavao, KP Bennett, « Évaluation de la confidentialité et de la qualité des soins de santé synthétiques Données » dans les actes de la Conférence sur l'intelligence artificielle pour la découverte et la réutilisation des données (ACM, New York, NY, États-Unis, 2019) ; <https://doi.org/10.1145/3359115.3359124>.
- 1344 X. Tang, R. Shin, HA Inan, A. Manoel, F. Mireshghallah, Z. Lin, S. Gopi, J. Kulkarni, R. Sim, « Privacy-Preserving In- « Apprentissage contextuel avec génération différentiellement privée de quelques plans » dans la 12e Conférence internationale sur les représentations d'apprentissage (2024) ; <https://openreview.net/forum?id=oZttOpRnOI>.
- 1345 F. Mireshghallah, Y. Su, T. Hashimoto, J. Eisner, R. Shin, « Adaptation de domaine préservant la confidentialité des analyseurs sémantiques » dans *ACL (1)* (2023), pp. 4950–4970 ; <https://doi.org/10.18653/v1/2023.acl-long.271>.
- 1346 J. Mattern, Z. Jin, B. Weggenmann, B. Schölkopf, M. Sachan, « Modèles de langage différentiellement privés pour des « Partage de données » dans *EMNLP* (2022), pp. 4860–4873 ; <https://aclanthology.org/2022.emnlp-main.323>.
- 1347 T. Stadler, B. Oprisanu, C. Troncoso, « Données synthétiques – Anonymisation Groundhog Day » dans 31e Symposium sur la sécurité USENIX (USENIX Security 22) (USENIX Association, Boston, MA, États-Unis, 2022), pp. 1451–1468 ; <https://www.usenix.org/conference/usenixsecurity22/presentation/stadler>.
- 1348 M. Meeus, F. Guepin, A.-M. Cretu, Y.-A. de Montjoye, « Talons d'Achille : identification des dossiers vulnérables dans « Publication de données synthétiques » dans le 28e Symposium européen sur la recherche en sécurité informatique (ESORICS 2023), G. Tsudik, M. Conti, K. Liang, G. Smaragdakis, éd. (Springer Nature Suisse, La Haye, Pays-Bas, 2024), pp. 380–399 ; https://doi.org/10.1007/978-3-031-51476-0_19.
- 1349 G. Ganey, E. De Cristofaro, Sur l'insuffisance des mesures de confidentialité basées sur la similarité : attaques de reconstruction contre « Données synthétiques véritablement anonymes », *arXiv [cs.CR]* (2023) ; <http://arxiv.org/abs/2312.05114>.

- 1350 R. Gilad-Bachrach, N. Dowlin, K. Laine, K. Lauter, M. Naehrig, J. Wernsing, « CryptoNets : application des réseaux neuronaux » dans les actes de la 33e Conférence internationale sur l'apprentissage automatique, MF Balcan, KQ Weinberger, éd. (PMLR, New York, New York, États-Unis, 2016) vol. 48 des Actes de la recherche sur l'apprentissage automatique, pp. 201–210 ; <https://proceedings.mlr.press/v48/gilad-bachrach16.html> .
- 1351 D. Kang, T. Hashimoto, I. Stoica, Y. Sun, « Mise à l'échelle de l'inférence DNN sans confiance avec des preuves à connaissance nulle » dans Atelier NeurIPS 2023 sur le ML régulier (La Nouvelle-Orléans, LA, États-Unis, 2023) ; <https://openreview.net/forum?id=GjNRF5VTfn>.
- 1352 B. Knott, S. Venkataraman, A. Hannun, S. Sengupta, M. Ibrahim, « CrypTen : le calcul multipartite sécurisé rencontre l'apprentissage automatique » dans Advances in Neural Information Processing Systems (Curran Associates, Inc., 2021) vol. 34, pp. 4961–4973 ; <https://papers.neurips.cc/paper/2021/hash/2754518221cfbc8d25c13a06a4cb8421-Résumé.html>.
- 1353 P. Mohassel, Y. Zhang, « SecureML : un système d'apprentissage automatique évolutif préservant la confidentialité » dans Symposium IEEE 2017 sur la sécurité et la confidentialité (SP) (IEEE Computer Society, San Jose, CA, États-Unis, 2017), pp. 19–38 ; <https://doi.org/10.1109/SP.2017.12>.
- 1354 O. Ohrimenko, F. Schuster, C. Fournet, A. Mehta, S. Nowozin, K. Vaswani, M. Costa, « Apprentissage automatique multipartite inconscient sur des processeurs de confiance » dans les actes du 25e symposium de la conférence USENIX sur la sécurité (SEC'16) (USENIX Association, Austin, TX, 2016), pp. 619–636 ; <https://www.usenix.org/conference/usenixsecurity16/technical-sessions/presentation/ohrimenko>.
- 1355 F. Tramer, D. Boneh, « Slalom : exécution rapide, vérifiable et privée de réseaux neuronaux dans du matériel de confiance » dans Conférence internationale sur les représentations d'apprentissage (2019) ; <https://openreview.net/forum?id=rJVorjCckKQ>.
- 1356* J. Zhu, H. Yin, P. Deng, A. Almeida, S. Zhou, Informatique confidentielle sur GPU nVIDIA H100 : une étude de performance Étude de référence, arXiv [cs.DC] (2024) ; <http://arxiv.org/abs/2409.03992>.
- 1357 T. South, J. Drean, A. Singh, G. Zyskind, R. Mahari, V. Sharma, P. Vepakomma, L. Kagal, S. Devadas, A. Pentland, « A. Feuille de route pour la confidentialité et la sécurité de bout en bout dans l'IA générative » (MIT, 2024) ; <https://doi.org/10.21428/e4baedd9.9af67664>.
- 1358 A. Cavoukian, Privacy by Design : Les 7 principes fondamentaux. (2009); <https://privacy.ucsc.edu/resources/privacy-by-design---foundational-principles.pdf>.
- 1359 M. ElBaih, Le rôle des réglementations sur la confidentialité dans le développement de l'IA (Une discussion sur les façons dont la confidentialité Les réglementations peuvent façonner le développement de l'IA) (2023) ; <https://doi.org/10.2139/ssrn.4589207>.
- 1360 E. Rader, R. Wash, B. Brooks, « Les histoires comme leçons informelles sur la sécurité » dans Actes du huitième symposium sur la confidentialité et la sécurité utilisables (ACM, New York, NY, États-Unis, 2012) ; <https://doi.org/10.1145/2335356.2335364>.
- 1361* J. Lamb, IA générative dans les soins de santé : tendances d'adoption et prochaines étapes (2024) ; <https://www.mckinsey.com/industries/healthcare/our-insights/generative-ai-in-healthcare-adoption-trends-and-whats-next> .
- 1362 G. Dhanuskodi, S. Guha, V. Krishnan, A. Manjunatha, M. O'Connor, R. Nertney, P. Rogers, Création des premiers GPU confidentiels : l'équipe de NVIDIA apporte confidentialité et intégrité au code et aux données utilisateur pour un calcul accéléré. Systèmes de mise en file d'attente. Théorie et applications 21, 68–93 (2023) ; <https://doi.org/10.1145/3623393.3623391>.
- 1363 X. Zhou, H. Kim, F. Brahma, L. Jiang, H. Zhu, X. Lu, F. Xu, BY Lin, Y. Choi, N. Mireshghallah, RL Bras, M. Sap, HAICOSYSTEM : un écosystème pour la gestion des risques de sécurité dans les interactions homme-IA, arXiv [cs.AI] (2024) ; <http://arxiv.org/abs/2409.16427>.
- 1364 K. Tirumala, AH Markosyan, L. Zettlemoyer, A. Aghajanyan, « Mémorisation sans surapprentissage : analyse de la dynamique d'apprentissage des grands modèles linguistiques » dans 36e Conférence internationale sur les systèmes de traitement de l'information neuronale (NeurIPS 2022) (Curran Associates Inc., Red Hook, NY, États-Unis, 2024) ; https://proceedings.neurips.cc/paper_files/paper/2022/file/fa0509f4dab6807e2cb465715bf2d249-Paper-Conference.pdf .
- 1365 N. Mireshghallah, H. Kim, X. Zhou, Y. Tsvetkov, M. Sap, R. Shokri, Y. Choi, « Les LLM peuvent-ils garder un secret ? Test des implications des modèles linguistiques sur la confidentialité via la théorie de l'intégrité contextuelle » dans ICLR (2024) ; <https://openreview.net/forum?id=gmg7t8b4s0>.
- 1366 M. Brundage, S. Avin, J. Wang, H. Belfield, G. Krueger, G. Hadfield, H. Khlaaf, J. Yang, H. Toner, R. Fong, T. Maharaj, P. W. Koh, S. Hooker, J. Leung, A. Trask, E. Bluemke, J. Lebensold, ... M. Anderljung, Vers un développement d'IA digne de confiance : mécanismes de soutien aux déclarations vérifiables, arXiv [cs.CY] (2020) ; <http://arxiv.org/abs/2004.07213>.



Toute demande de renseignements concernant cette publication doit être envoyée à : Secretariat.aistateofscience@dsit.gov.uk

Numéro de série de recherche : DSIT 2024/000

Publié en janvier 2025 par le gouvernement britannique

@Droits d'auteur de la Couronne 2025