

---

L'IA : une course effrénée à la donnée

## Description

À la fin de l'année 2021, OpenAI, à l'origine de ChatGPT, a été confronté à un problème de taille pour entraîner la nouvelle version de son intelligence artificielle (IA) générative, GPT-4 : l'épuisement des contenus en anglais de qualité, l'intégralité ayant déjà été utilisée.

La conclusion d'une enquête menée par cinq journalistes du *New York Times* parue en avril 2024 est sans appel : « OpenAI, Google et Meta ont ignoré les politiques d'entreprise, altéré leurs propres règles et ont discuté de la possibilité de contourner la loi sur le droit d'auteur alors qu'ils cherchaient des informations en ligne pour entraîner leurs derniers systèmes d'intelligence artificielle. »

En janvier 2020, Jared Kaplan, physicien théoricien à l'université Johns-Hopkins, a publié un article majeur dans le domaine de l'IA qui démontre qu'un modèle neuronal est caractérisé par quatre paramètres : la taille du modèle, la taille du jeu de données, l'entraînement, le coût de l'entraînement et la performance après l'entraînement, et que ces paramètres sont liés par des lois statistiques dont l'une stipule que plus on dispose de données pour entraîner un grand modèle de langage, mieux il fonctionnera. « Tout le monde a été très surpris que ces tendances – ces lois d'échelle comme nous les appelons – soient essentiellement aussi précises que ce que l'on observe en astronomie ou en physique », a-t-il déclaré après la publication de cet article, co-crit avec neuf chercheurs d'OpenAI.

Depuis, les géants de la tech impliqués dans l'IA sont engagés dans une course effrénée à la donnée. Les grands modèles de langage sont en effet entraînés à partir d'immenses corpus de données textuelles, contenant plusieurs milliards de mots, représentés sous la forme de tokens, afin d'être manipulables par le programme informatique ([voir La rem n°65-66, p.27](#)). Selon Pablo Villalobos, de l'Institut de recherche Epoch, cité par le *Wall Street Journal*, ChatGPT-4 a été entraîné à partir de 12 000 milliards de tokens et, si ChatGPT-5 suit la même évolution que le passage de ChatGPT-3 à ChatGPT-4, 60 000 à 100 000 milliards de tokens seront nécessaires. Or, il n'existe plus assez de nouvelles données de qualité pour entraîner ces IA, c'est-à-dire des données issues de livres, d'articles de presse, de publications scientifiques, de code informatique open source, de contenus provenant de Wikipédia ainsi que de tous les corpus de textes spécialement dédiés à l'entraînement de ces programmes.

En 2022, des chercheurs de l'université d'Aberdeen en Écosse, de l'université de

Tübingen en Allemagne et du MIT (Massachusetts Institute of Technology) aux États-Unis avaient d'ailleurs prouvé ce phénomène de rarefaction des contenus à la fois de qualité et non protégés par le droit d'auteur ou par les réglementations en vigueur sur les données personnelles ([voir La rem n°65-66, p.27](#)).

### Plusieurs solutions existent alors aux développeurs de ces IA

Ils peuvent passer des contrats avec des fournisseurs de contenus, en particulier les médias ; entraîner les IA sur des contenus créés par des IA ou bien utiliser sans autorisation des contenus protégés par le droit d'auteur. Alors que, depuis décembre 2023, l'entreprise fait toujours l'objet de poursuites pour violation des droits d'auteur de la part du *New York Times*, OpenAI a conclu un premier accord pluriannuel avec l'agence Associated Press la même année. En 2024, des accords similaires ont été signés avec *Le Monde* en France, Prisa pour *El País* en Espagne, Springer pour le *Bild*, *Die Welt*, *Politico* en Allemagne, et récemment avec le *Financial Times* en Angleterre. Si peu d'informations filtrent sur ces accords ; le *Financial Times* avait justement révélé, concernant le groupe allemand, la signature d'un contrat se chiffrant en dizaines de millions d'euros par an.

Une autre solution pour pallier la pénurie de contenus serait d'entraîner des IA sur des données synthétiques, c'est-à-dire des données elles-mêmes produites par une IA. « Nous avons remarqué que si le modèle est suffisamment bon et que la quantité de données générées n'est pas trop importante par rapport aux réelles, alors le modèle ne dérangera pas », explique Quentin Bertrand, chercheur au Mila et à l'Université de Montréal, coauteur d'une étude publiée en janvier 2024 et intitulée « Stabilité du recyclage itératif des modèles génératifs sur leurs propres données ». Même si le risque d'un effondrement du modèle, *model collapse*, est réel, et que certains chercheurs considèrent les données synthétiques comme l'équivalent informatique de la consanguinité, toutes les grandes entreprises du secteur travaillent activement sur le sujet.

Autre option, enfin : se servir, sans autorisation, de contenus protégés par le droit d'auteur. En 2021, les chercheurs d'OpenAI ont créé un outil de reconnaissance vocale, appelé Whisper, capable de transcrire l'audio des vidéos hébergées sur YouTube sous la forme de textes, afin de mettre la main sur des contenus inédits. Bien évidemment, Google interdit tout autant l'utilisation par des tiers des vidéos publiées sur la plateforme que l'usage d'outils permettant un accès automatisé à ces mêmes vidéos. Or, selon l'enquête publiée par le *New York Times*, « Certains employés de Google savaient qu'OpenAI avait exploité des vidéos de YouTube pour obtenir des données [à] Mais ils n'ont pas arrêté OpenAI parce que Google avait également utilisé des transcriptions de vidéos YouTube pour entraîner ses modèles d'intelligence artificielle ». Alors que lui-même violait les droits d'auteur des utilisateurs de sa propre plateforme, Google ne pouvait prendre le risque d'accuser OpenAI de faire la même chose.

Open AI, Meta ou encore Google sont donc confrontés en même temps au même problème : la pénurie de données de qualité pour entraîner leur IA. *« La seule chose qui nous empêche d'être aussi bons que ChatGPT, c'est littéralement le volume de données »*, déclarait Nick Grudin, vice-président des partenariats et du contenu mondial de Meta, dans des propos rapportés par le *New York Times*. Des enregistrements de réunions internes à Meta datant de 2023, impliquant des responsables, des avocats et des ingénieurs obtenus par *The Times*, indiquent également des discussions au sujet de la collecte de données protégées par le droit d'auteur sur l'ensemble du web en assumant le risque de poursuites judiciaires, la négociation de licences avec les éditeurs, les artistes, les musiciens et l'industrie de l'information tant jugés trop chronophages.

OpenAI reconnaît d'ailleurs ouvertement l'utilisation de documents protégés par le droit d'auteur. Dans un témoignage écrit, fourni dans le cadre d'une enquête sur les grands modèles de langage menée par la Commission des communications et du numérique de la Chambre des Lords du Parlement anglais, l'entreprise tient ces propos : *« Notamment que le droit d'auteur couvre aujourd'hui pratiquement toutes les formes d'expression humaine »* y compris les articles de blog, les photographies, les messages de forum, les bouts de code de logiciel et les documents gouvernementaux, *« il serait impossible d'entraîner les principaux modèles d'IA actuels sans utiliser des documents protégés par le droit d'auteur »*.

Le principal argument juridique, systématiquement invoqué par OpenAI pour justifier ces violations du droit d'auteur, repose sur le *fair use*. En décembre 2023, en réponse à la plainte pour violation du copyright déposée par le *New York Times*, OpenAI considère qu'*« il est clair depuis longtemps que l'utilisation non-consommatrice de matériel protégé par le copyright (comme la formation de grands modèles de langage) est protégée par le fair use »*. Ce principe juridique de propriété intellectuelle propre au droit américain d'un usage loyal ou équitable autorise, dans des circonstances précises, à utiliser des œuvres protégées par le droit d'auteur sans l'autorisation préalable du titulaire des droits, à condition que cet usage soit considéré comme raisonnable et équitable. Au-delà du raisonnement juridique, qui semble tout à fait inapproprié si l'on considère les millions de livres ou d'articles et les millions d'heures de vidéos YouTube protégés par le droit d'auteur pris sans autorisation par ces IA et l'activité dorénavant purement commerciale de l'entreprise, l'argumentation d'OpenAI consiste également à minimiser l'importance et le volume des contenus collectés sans autorisation auprès d'une entreprise en les rapportant à l'intégralité des données sur lesquelles le modèle d'IA est entraîné.

Autrement dit, puisque la ligne droite est plus que la somme des points qui la composent, pour reprendre une idée d'Aristote, le contenu du *New York Times*, comme toute source unique, ne *« contribue pas de manière significative »* à l'entraînement de ChatGPT, aurait déclaré OpenAI au média américain. Et une fois qu'un contenu a servi à entraîner une IA, il est extrêmement complexe de le « retirer » ou de le faire « désapprendre » ou même de montrer ou

---

prouver quoi que ce soit puisque ces programmes sont des boîtes noires, y compris pour leurs concepteurs ([voir La rem n°68, p.41](#)). On se demande de quelle loyauté ou équité se prévalent les dirigeants de ces géants de la tech qui, sans scrupule, foulent aux pieds les législations nationales et internationales, au motif que celui qui aura fait avaler le plus de données à son IA remportera le marché.

Sources :

- UK Parliament, « OpenAI » written evidence (LLM0113). House of Lords Communications and Digital Select Committee inquiry: Large language models », December 5, 2023.
- Thomas Daniel, Murgia Madhumita, « Axel Springer strikes landmark deal with OpenAI over access to news titles », ft.com, December 13, 2023.
- Marin Joréme, « Pour éviter de nouveaux procès, OpenAI négocie avec des éditeurs de presse », usine-digitale.fr, 5 janvier 2024.
- Saramour Céilia, « Les œuvres soumises au droit d'auteur indispensables pour entraîner ChatGPT, admet OpenAI », usine-digitale.fr, 9 janvier 2024.
- Clavey Martin, « OpenAI contre-attaque et accuse le *New York Times* d'avoir hacké ses produits », next.ink, 29 février 2024.
- Metz Cade, Kang Cecilia, Frenkel Sheera, Thompson Stuart A., Grant Nico, « How Tech Giants Cut Corners to Harvest Data for A.I. », *The New York Times*, nytimes.com, April 6, 2024.
- Seetharaman Deepa, « For Data-Guzzling AI Companies, the Internet Is Too Small », wsj.com, April 1, 2024.
- Saramour Céilia, « Après Axel Springer et *Le Monde*, OpenAI s'achète les faveurs du *Financial Times* », usine-digitale.fr, 30 avril 2024.

## Categorie

1. Techniques

**date création**

4 juillet 2024

**Auteur**

jacquesandrefines