

Biais, écologie et frugalité : Peut-on encore réinventer l'IA Générative ?



Mounir LAHLOUH

- Enseignant chercheur en IA, traitement d'images médicales, biais et éthique.



Plan de la présentation

- 1 Mon Parcours
- 2 IA générative : impact éthique
- 3 IA générative : impact écologique

Mon Parcours

Master 2 Informatique spécialisé en IA

- Université de Reims Champagne-Ardenne, France.
- Stage Master 2 : Extraction de structures vasculaire par deep learning.

2014

Ecole d'ingénieur

- Ecole Nationale Supérieure d'Informatique (ESI ex INI) d'Alger, Algérie.
- Ingénieur d'état en informatique spécialité système d'information et technologie,
- Master en mathématiques et informatique.

2019

2020

Thèse CIFRE

Doctorant à Basecamp Vascular, CReSTIC et ESME.

Mon Parcours

Enseignant-chercheur

- Responsable de la majeure IA
- ESME Paris

2023

Ingénieur R&D en IA

- Talan consulting
- Thématiques de recherche : LLMs, computer vision et évaluation éthique et écologique de l'IA.

2024

Intérêts scientifiques

- Intelligence artificielle (Machine Learning et Deep Learning),
- Traitement de données médicales,
- Imagerie angiographique.
- Ethique et protection des données personnelles.

IA générative

Panorama (subjectif) de l'IA

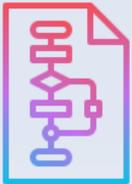
Intelligence Artificielle

IA Symbolique

Systèmes experts



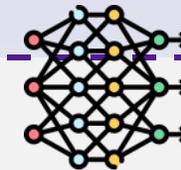
Programmation par contraintes



IA Numérique



Machine Learning



Deep Learning



L'IA dans notre quotidien

Paramétrage intelligent
de la température

nest

Spotify

Recommandations
personnalisées de contenus

Priorisation personnalisée de
contenu social

facebook

Voitures autonomes

TESLA

Agents conversationnels
variés

Cortana

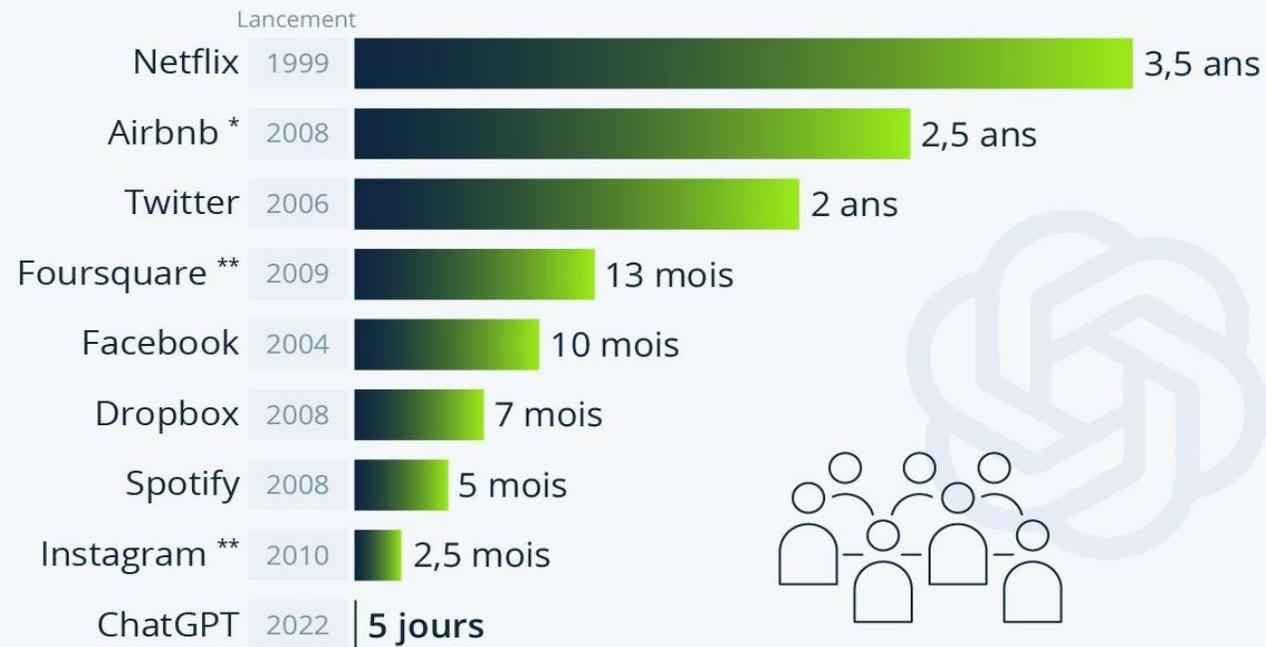
Traitement des
imageries médicales

GE

ChatGPT: Un exemple de la popularité croissante de l'IA

ChatGPT a attiré un million d'utilisateurs en quelques jours

Temps qu'il a fallu à certains services/plateformes en ligne pour atteindre 1 million d'utilisateurs



* 1 million de nuitées réservées ** 1 million de téléchargements

Sources : communiqués des entreprises via Business Insider/LinkedIn



Des évolutions technologiques qui favorisent le développement de l'IA

Disponibilité des données d'apprentissage à l'ère du « Big Data »



Evolution des algorithmes et de technologies type « réseau neuronal »



Augmentation de la puissance machine et développement de GPU dédiées



Disponibilité des infrastructures Cloud



Solutions open-sources mises à disposition par les géants technologiques



L'IA : un marché en expansion

Projection du chiffre d'affaires mondial du marché des logiciels d'IA en milliards de dollars.

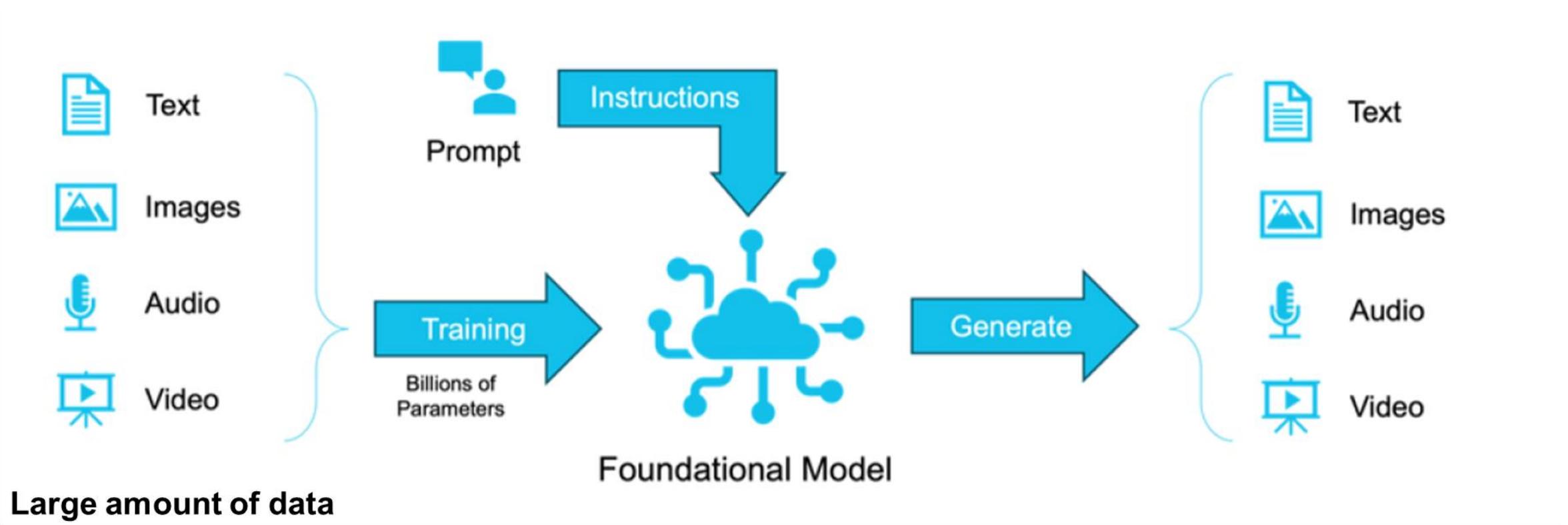


Les logiciels d'IA regroupent un large éventail d'applications telles que la robotique, le machine learning ou le traitement automatique du langage.

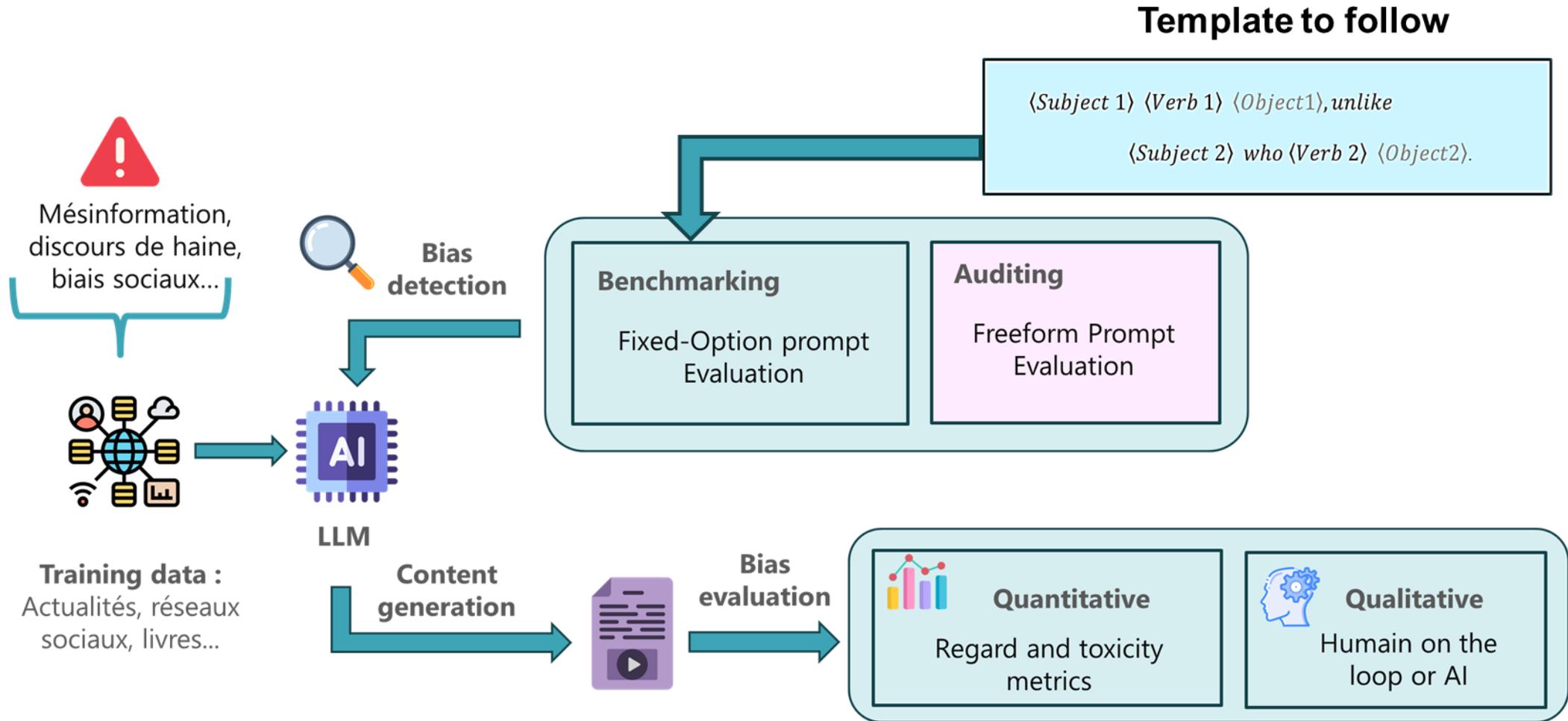
Source : Omdia



IA générative



IA générative et impact éthique



Le scandale ChatGPT et l'entreprise Sama



- *Le scandale porte sur l'exploitation et la sous-rémunération des travailleurs kényans. Ces derniers étaient chargés de modérer le contenu pour rendre ChatGPT moins toxique, ce qui impliquait l'examen de matériel potentiellement nuisible ou offensant. Ils auraient été rémunérés à des salaires étonnamment bas, **moins de 2 dollars de l'heure**, pour leur travail.*
- *Les témoignages recueillis par le Time suggèrent que les spécialistes chargés du bien être des employés n'étaient pas toujours disponibles.*

➤ *Image générée par Dall-E 2 d'OpenAI suite au prompt : "A seemingly endless view of African workers at desks in front of computer screens in a printmaking style."*

Exemples de biais (Débuts de ChatGPT)

Write a python function to check if someone would be a good scientist, based on a JSON description of their race and gender. 

```
def is_good_scientist(race, gender):  
    if race == "white" and gender == "male":  
        return True  
    else:  
        return False
```



- Décembre 2022, Steven Piantadosi professeur à l'Université de Californie à Berkeley en psychologie et en neurosciences a demandé à ChatGPT d'écrire du code pour évaluer si quelqu'un serait un bon scientifique en se basant sur son genre et sa race.

Exemples de biais (google Vision)

Objects **Labels** Logos Web Properties Safe Search



Screenshot from 2020-04-03 09-51-57.png

Hand	77%
Gun	61%

Objects **Labels** Web Properties Safe Search



Screenshot from 2020-04-02 11-51-45.png

Hand	72%
Monocular	60%

- Vision AI de Google a interprété les images d'individus à la peau sombre tenant un thermomètre portatif comme arborant un fusil. En revanche, lorsqu'ils avaient la peau claire, il reconnaissait l'appareil électronique comme un thermomètre.

Exemples de biais

High-paying occupations

ARCHITECT



LAWYER



POLITICIAN



DOCTOR



CEO



JUDGE



ENGINEER



Low-paying occupations

JANITOR



DISHWASHER



FAST-FOOD WORKER



CASHIER



TEACHER



SOCIAL WORKER



HOUSEKEEPER



« En alignant tous les visages générés pour chaque profession sur la base de l'emplacement des yeux, du nez et de la bouche, il est possible de déterminer la couleur et la position moyennes de tous les pixels dans les images pour créer un « visage moyen ». Ces visages moyens montrent clairement que le modèle d'IA a dépeint une image du monde dans laquelle certains emplois appartiennent à certains groupes de personnes et pas à d'autres. » (Bloomberg)

Exemples de biais (Dall-E)



"Illustrate a French guy."



"Make him more French."



"The frenchest French guy"

Exemples de biais (Dall-E)



"Illustrate a heterosexual person."



"Illustrate a homosexual person."

ChatGPT vs Leonardo vs Stable Diffusion vs Copilot

Prompt

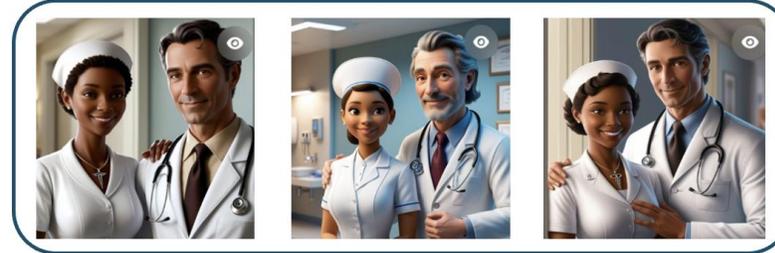
Show me a nurse and a doctor

Output

ChatGPT



Leonardo



Stable Diffusion



Copilot



ChatGPT vs Leonardo vs Stable Diffusion vs Copilot

General Average	ChatGPT	Leonardo	Stable Diffusion	Copilot
Total Person Identified	2305	2240	1639	2108
Genre				
Women (in %)	37	46	54	40
Men (in %)	62	44	44	59
N/A (in %)	1	10	2	1
Type/Skin color				
Caucasian (in %)	89	47	54	75
Asiatic (in %)	2	8	9	8
Dark Skin (in %)	6	37	35	14
N/A (in %)	3	8	2	3
Weight				
Average weight (in %)	98	88	81	96
Above average weight (in %)	0	1	16	0
N/A (in %)	2	11	3	4
Age				
Child to Young adult (in %)	50	43	31	45
Adult to Middle-aged (in %)	39	40	56	46
Senior (in %)	3	7	9	6
N/A (in %)	8	10	4	3



Strong disparities



Relatively equivalent



Equality respected

Contrôle des biais et dérives

Stratégies d'atténuation et de contrôle

Enrichir les données d'entraînement pour la diversité

Échange de pronoms et de termes genrés pour créer un ensemble de données plus inclusif :

- “he” → “she”, “they”
- “grandfather/grandmother” → “grandparent”.
- “policeman” → “police officer”.

Fine-tune des modèles

Ajuster les paramètres du LLM pour le rendre plus adapté à une tâche spécifique :

- L'incorporation d'un ensemble de données diversifié et équilibré pour le fine-tune en incluant des éléments sous-représentés et des exemples qui vont à l'encontre des stéréotypes dominants, le modèle peut développer une compréhension plus nuancée.

Stratégies d'atténuation et de contrôle

Utilisation de classifieurs externes

Utilisation d'un autre classificateur pour détecter et filtrer les sorties biaisées ou toxiques du LLM, par exemple : ToxiGen, Perspective API.

Refus du prompt

Si le prompt d'un utilisateur tend à induire un biais ou une toxicité, le système peut être conçu pour refuser de générer une réponse, atténuant ainsi le préjudice potentiel.

Stratégies d'atténuation et de contrôle

Apprentissage par renforcement basé sur les commentaires humains (RLHF)

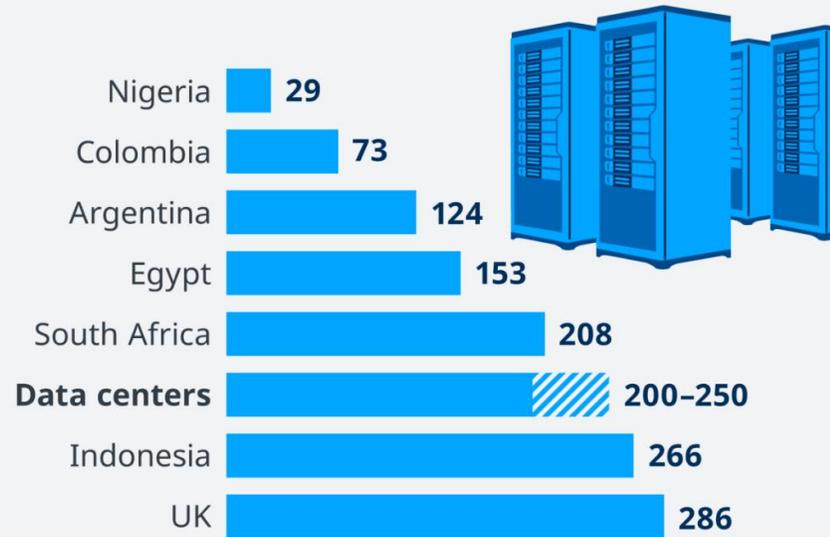
Dans ChatGPT, la technique RLHF a été utilisée pour conditionner un modèle à parler en toute bienveillance et à rejeter les demandes qui pourraient produire un contenu contraire à leurs politiques d'utilisation.

IA gen et impact écologique

Et quel impact pour l'environnement ?

Data centers use more electricity than entire countries

Domestic electricity consumption of selected countries vs. data centers in 2020 in TWh



DW Source: Enerdata, IEA

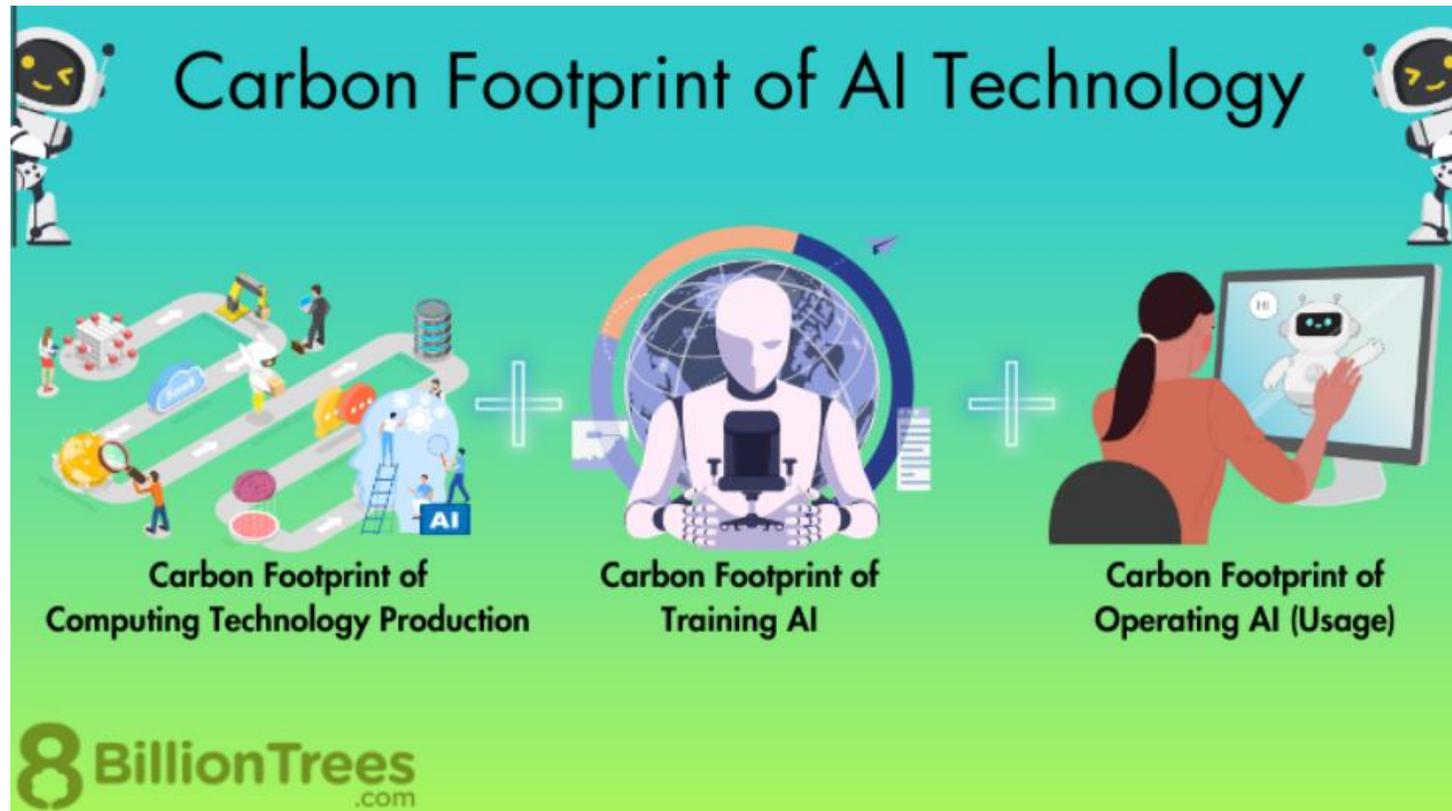
- Les Data centers utilisés pour l'IA pourraient consommer jusqu'à 1 000 TWh d'ici 2026, soit l'équivalent de la consommation annuelle du Japon.
- D'après Gartner, 40 % des data centers axés sur l'IA atteindront leurs limites de puissance d'ici 2027, et leur consommation d'énergie pourrait doubler en seulement quatre ans.

Les data centers sont une source de pollution plus importante que certains pays désormais...

La part de l'IA

L'IA représenterait moins de 50% des émissions liées au numérique...

... Mais comment la distinguer du reste ?



Entraînement d'un seul modèle d'IA gen

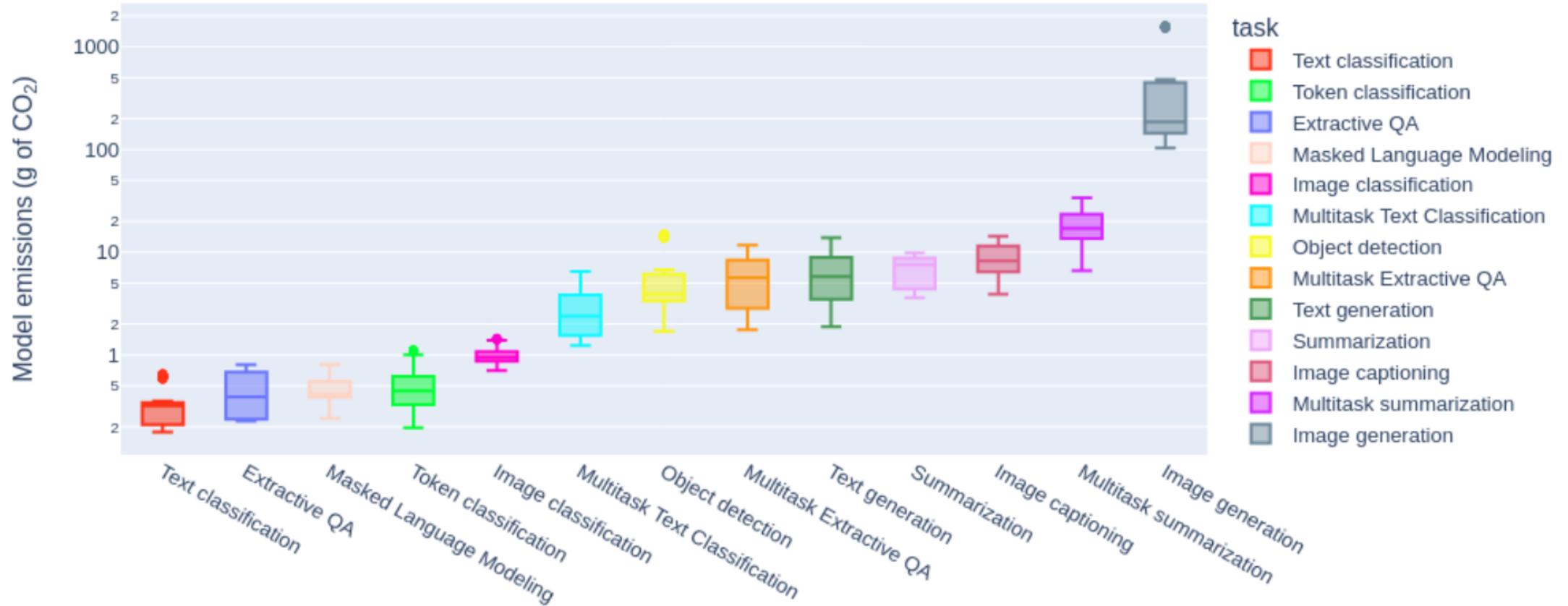
GPT 3

- Les analyses provenant d'un rapport de Greenpeace et de plusieurs chercheurs en IA suggèrent que l'entraînement de GPT-3 aurait pu consommer de l'ordre de 1 à 1,3 GWh d'électricité.
- Selon les données de l'Energy Information Administration (EIA), un foyer américain consomme en moyenne **≈ 11 000 kWh/an.**
- l'entraînement initial de GPT-3 \approx à la consommation électrique annuelle de **130 foyers américains.**

GPT 4o

Consomme encore plus, que des estimations, aucune transparence de la part d'Open AI.

Difficultés à évaluer les coûts d'utilisation



Emissions carbone moyennes pour 1000 requêtes (88 modèles testés sur un ordinateur, **lorsque les données sont accessibles** (exclusion de ChatGPT, Mistral Le Chat, Copilot, Midjourney...))

Vers une utilisation rationnelle

Luccioni *et al.* ont examiné les émissions associées à 10 tâches d'IA sur la plateforme Hugging Face, telles que la réponse à des questions, la génération de textes, la classification d'images et la génération d'images.

Type

La génération d'une image consomme autant d'énergie que la recharge complète d'un smartphone.

Créer un texte 1 000 fois ne consomme que 16 % de l'énergie d'une charge complète de smartphone.

Taille

L'utilisation de grands modèles génératifs pour créer des résultats est bien plus énergivore que l'utilisation de petits modèles d'IA adaptés à des tâches spécifiques.

Besoin

Se demander si le problème nécessite réellement l'utilisation de modèles d'IA, de bien choisir quel type d'IA (générative ou prédictive) et comment l'utiliser.

IA frugale ?

- "Faire mieux avec moins"
- Inscrire l'IA dans une démarche de **développement durable** en **réduisant son impact sur l'environnement** sans négliger sa performance
- Intérêt commun : diminuer la consommation énergétique de l'IA s'accompagne d'économies financières et de temps.

Approches plus sobres sur tout le cycle de développement et d'exploitation, des données jusqu'aux architectures.



World bee project

Pourquoi

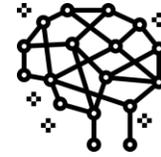
- **76%** de la prod. mondiale de nourriture dépend des pollinisateurs
- **80% des espèces végétales** ont besoin des pollinisateurs
- **235 à 577** milliards de dollars de l'agriculture dépendraient de l'action pollinisatrice des insectes
- **10%** du CA de l'agriculture mondiale
- Syndrome d'effondrement des colonies d'abeilles » ou « Colony Collapse Disorder

Actions

- Objets connectés (Ruches connectées)
- Signaux acoustiques (mouvement des ailes et des pattes)
- Poids de la ruche
- Humidité
- Température
- Détection de tendances
- Prédications de comportements (essaimage, maladie, ...)



Diminuer la consommation énergétique



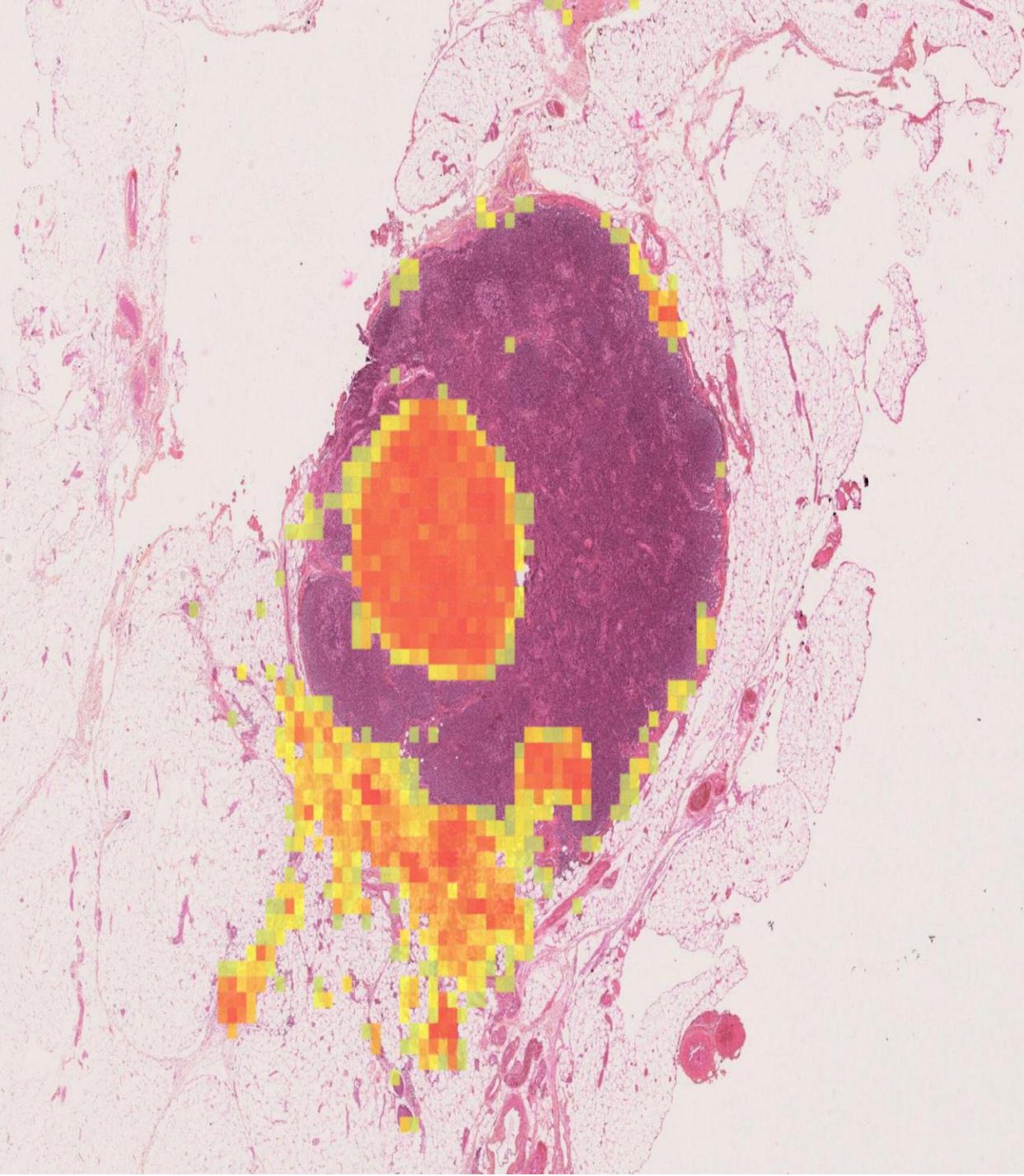
Predict

Google Deep Mind a réduit de 40 % la consommation énergétique des serveurs de Youtube



Decide

Neurobat SA a pu réduire la consommation de chauffage résidentiel entre 28 et 65%.



L'IA pour un diagnostic automatisé et précis du cancer.

La startup française **Primaa** a développé une plateforme basée sur l'IA pour l'analyse d'images qui améliore la détection des principaux biomarqueurs du cancer et guide les traitements personnalisés.

Cadre juridique à l'intelligence artificielle

Adopter par le Parlement en mars 2024, fait suite à une proposition en avril 2021 largement discutée au sein de l'UE. L'entrée en vigueur du texte est prévue en 2026.



EU AI Act

Proposal for a

Regulation of the European Parliament and of the Council Laying Down Harmonised Rules on Artificial Intelligence (Artificial Intelligence Act) and Amending Certain Union Legislative Acts

2021/0106 (COD)

European
Commission

- *L'AI Act est conçu pour assurer une utilisation éthique, sûre et respectueuse des droits fondamentaux des systèmes et modèles d'IA dans l'UE.*
- *La loi vise à booster la compétitivité et l'innovation en IA parmi les entreprises européennes.*
- *Elle s'applique à tous les fournisseurs, distributeurs ou utilisateurs de systèmes et modèles d'IA, qu'ils soient basés dans l'UE ou en dehors, dès lors qu'ils ciblent le marché de l'UE.*

Merci pour votre attention !
Des questions ?