

ENSEIGNER L'IA, ENSEIGNER AVEC L'IA : DEUX DÉFIS MAJEURS POUR AGROPARISTECH

Mardi 21 Octobre 2025
Journées d'Octobre de l'UPA

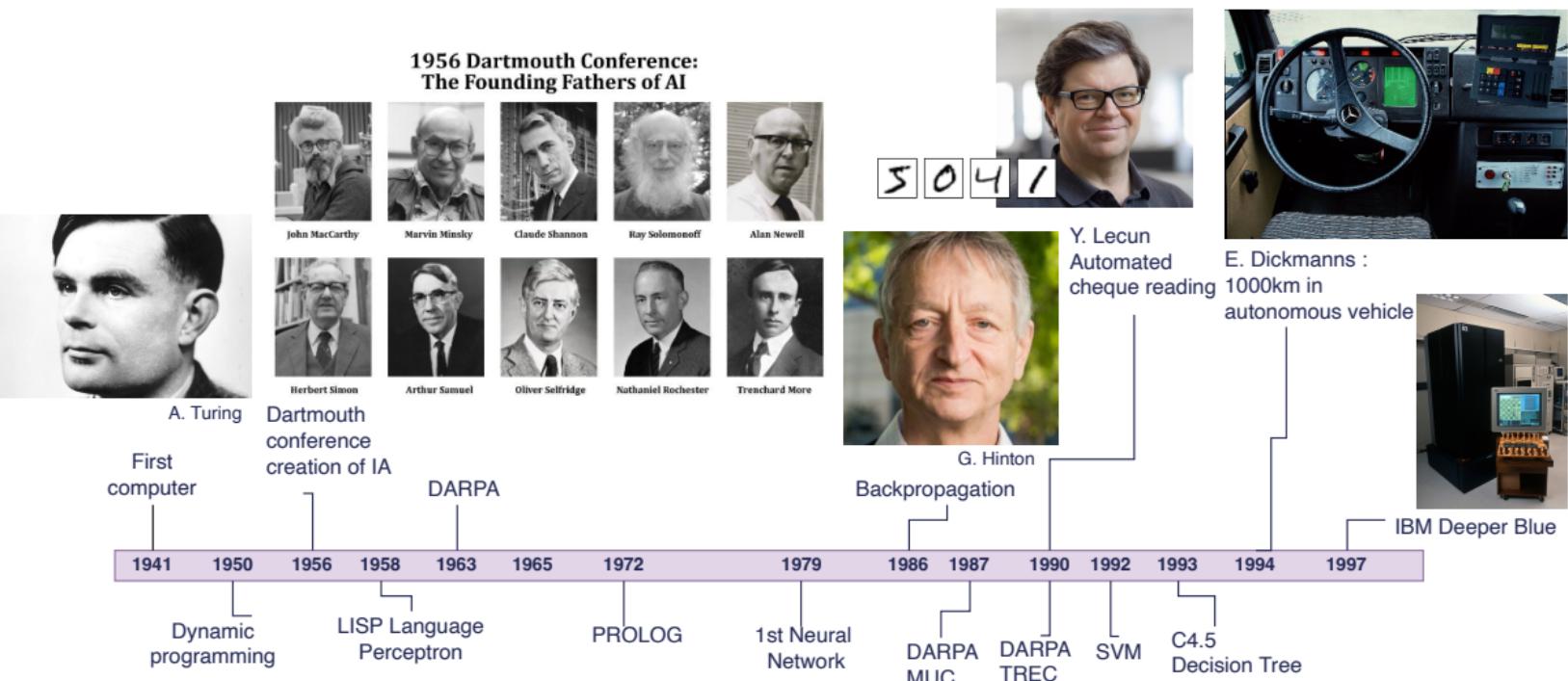
Vincent Guigue
<https://vguigue.github.io>

INTRODUCTION



Un rapide tour historique de l'Intelligence Artificielle

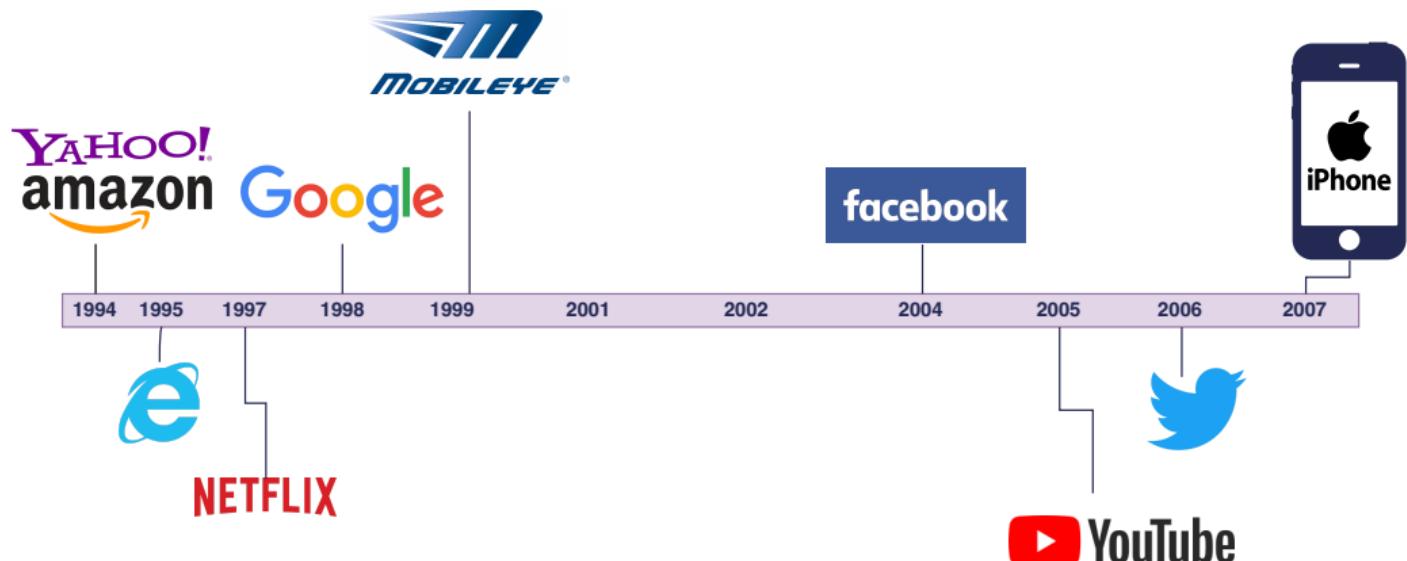
Naissance de l'informatique... Et de l'Intelligence Artificielle





Un rapide tour historique de l'Intelligence Artificielle

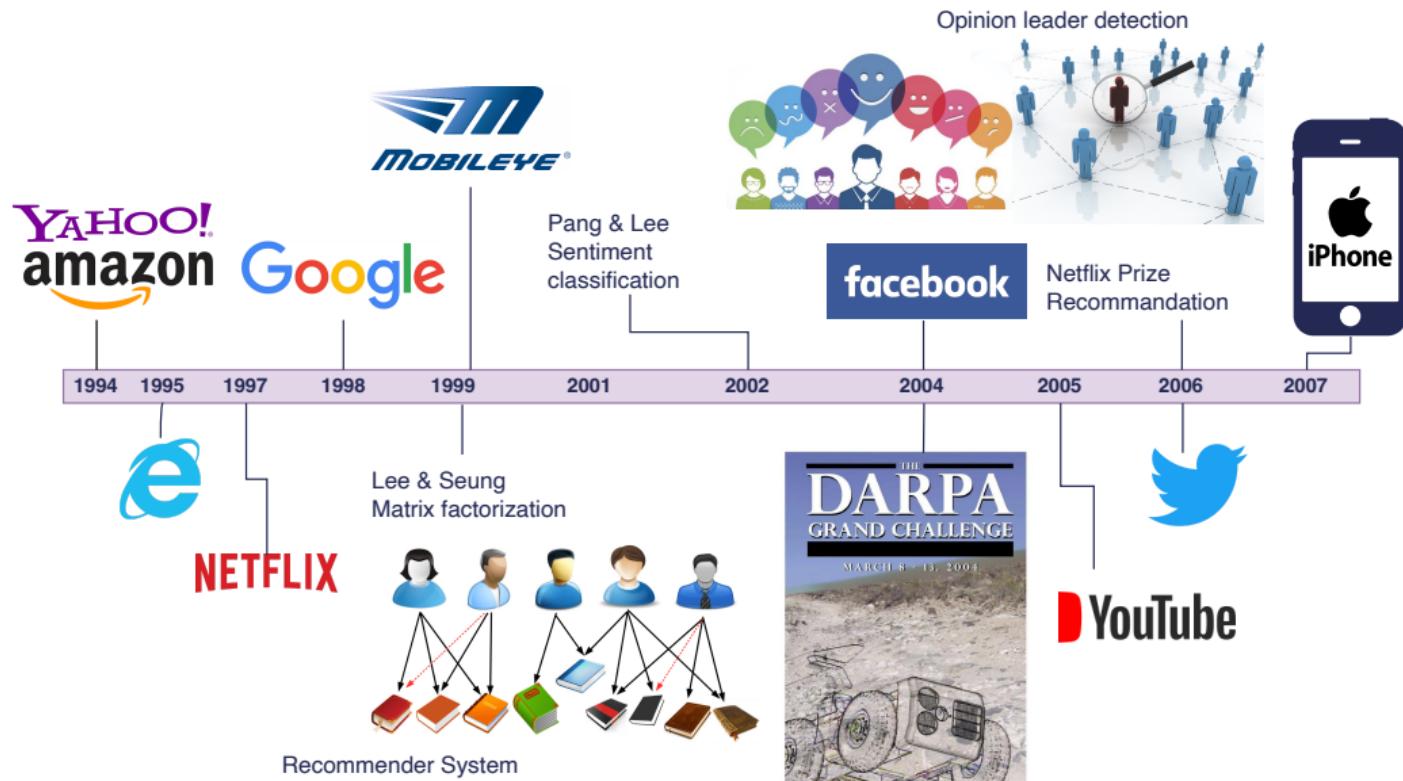
Emergence (ou refondation) des GAFAM/GAMMA





Un rapide tour historique de l'Intelligence Artificielle

Emergence (ou refondation) des GAFAM/GAMMA



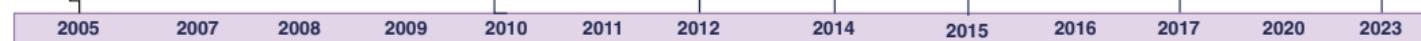


Un rapide tour historique de l'Intelligence Artificielle

Formation d'une vague de l'Intelligence Artificielle



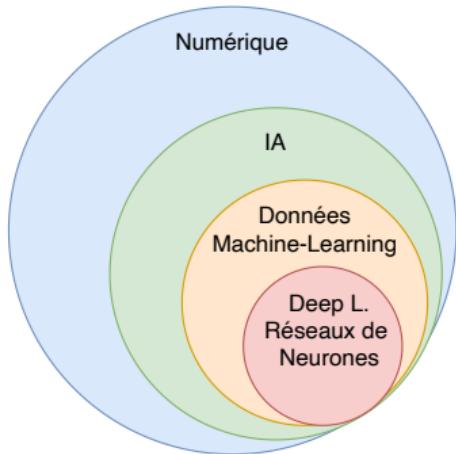
Thrun:
DARPA Gd Challenge
victory



IBM Jeopardy win



Artificial Intelligence & Machine Learning



Input (X)	Output (Y)	Application
email	spam? (0/1)	spam filtering
audio	text transcript	speech recognition
English	Chinese	machine translation
ad, user info	click? (0/1)	online advertising
image, radar info	position of other cars	self-driving car
image of phone	defect? (0/1)	visual inspection

IA : programmes informatiques qui s'adonnent à des tâches qui sont, pour l'instant, accomplies de façon plus satisfaisante par des êtres humains car elles demandent des processus mentaux de haut niveau.

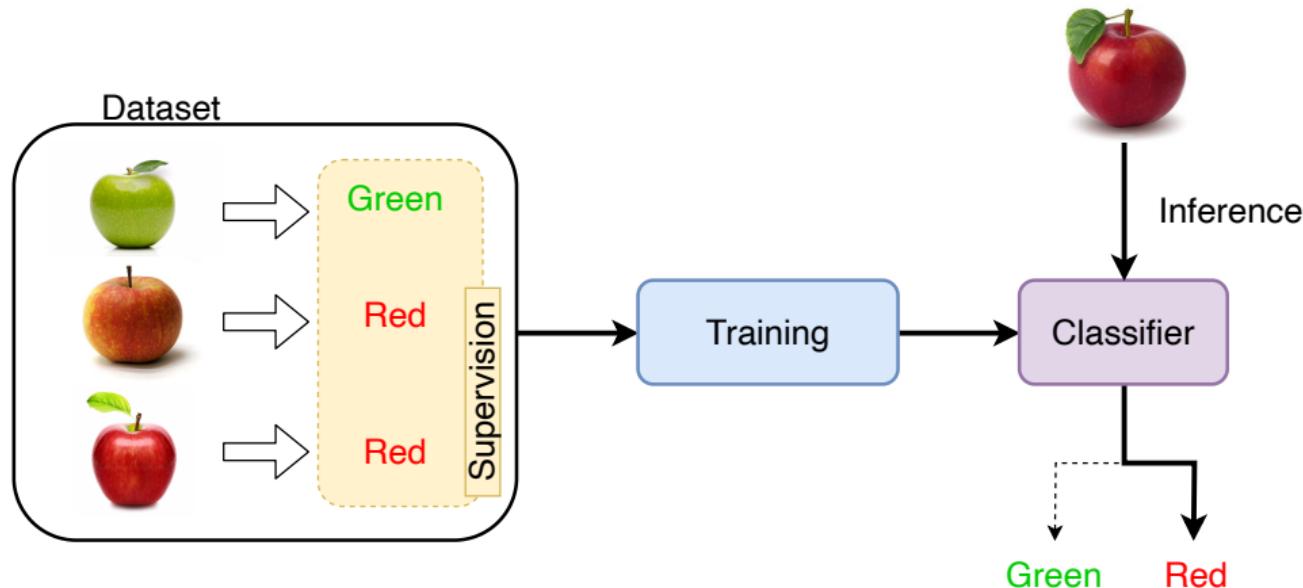
Marvin Lee Minsky, 1956

N-AI (Narrow Artificial Intelligence), dédiée à une tâche unique

≠ **IA-G (IA Générale)**, qui remplace l'humain dans les systèmes complexes.

Andrew Ng, 2015

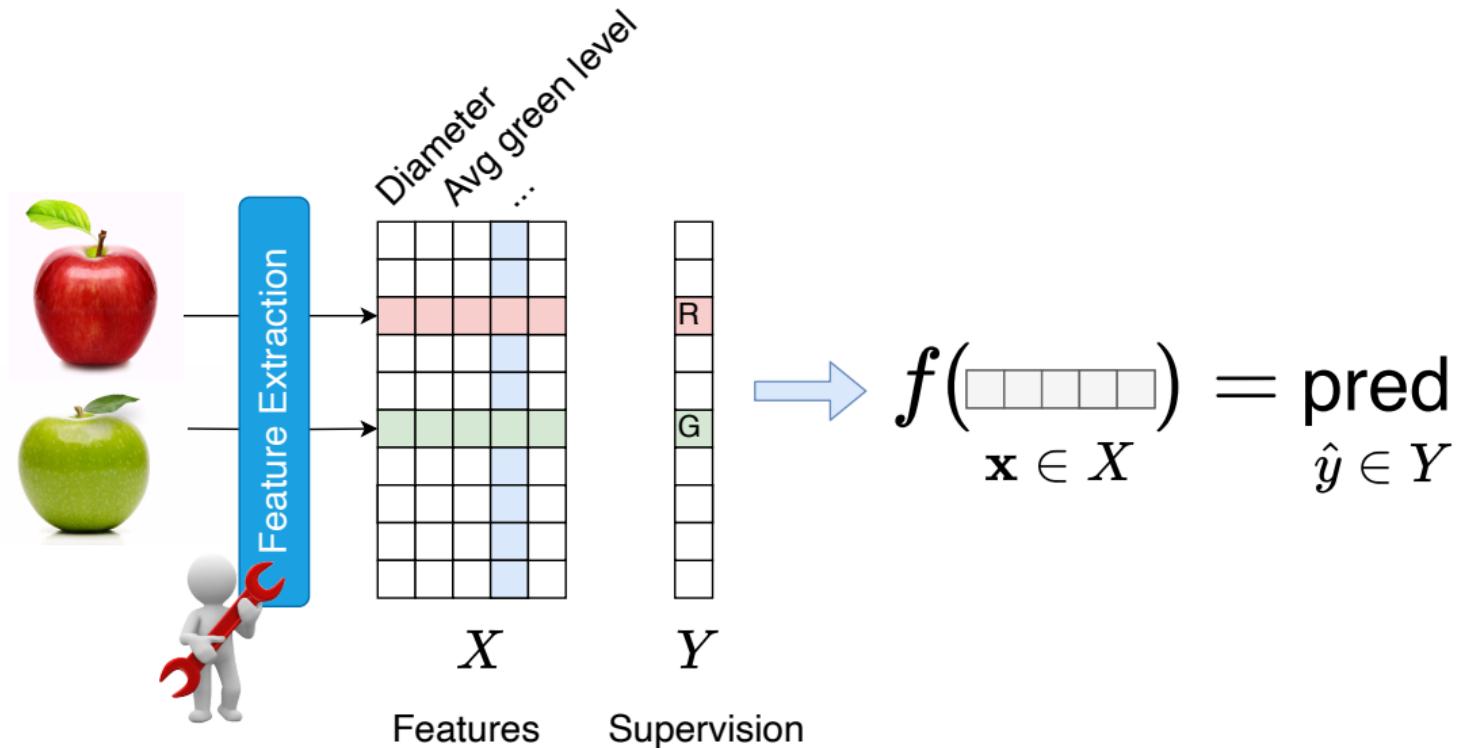
Chaîne de Traitement Supervisé & Modèles



- Promesse = construire un modèle *uniquement* à partir d'observations

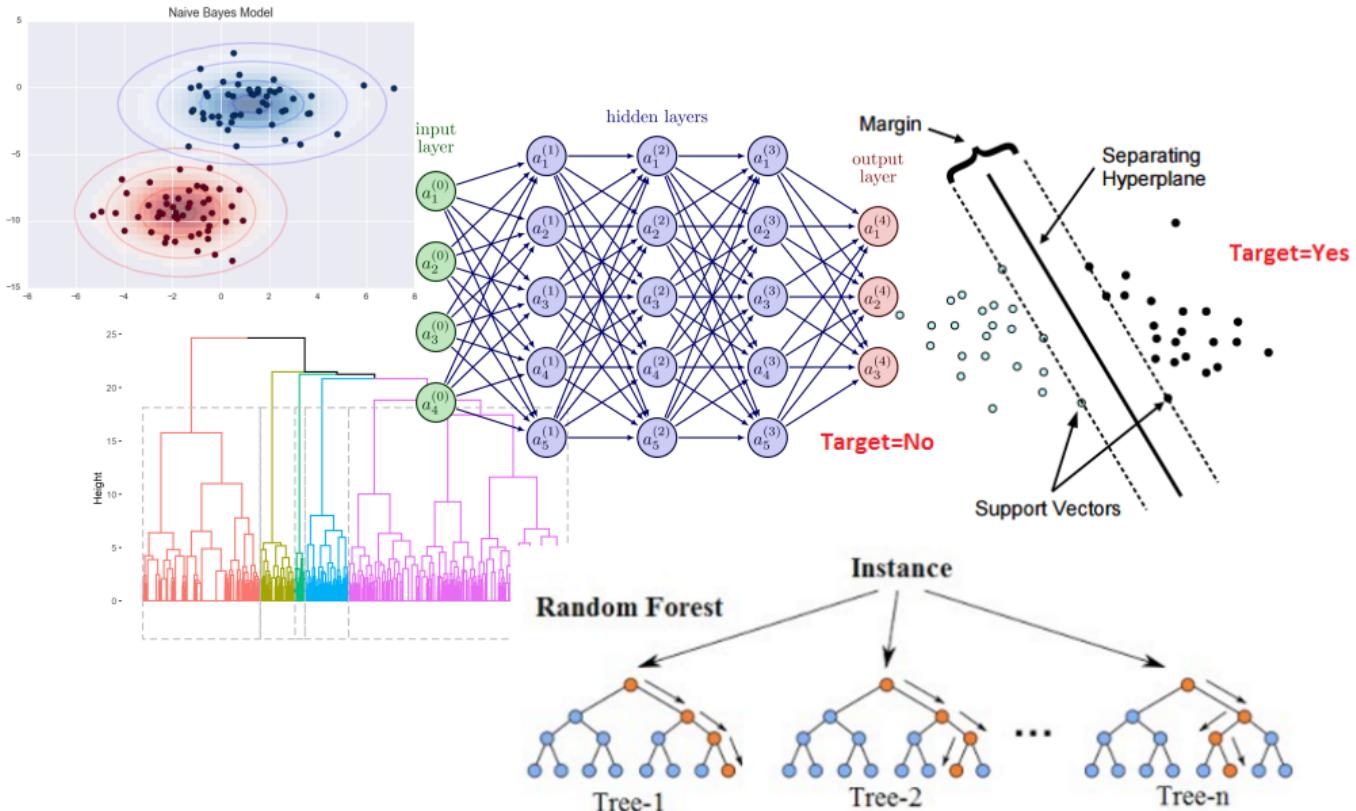


Chaîne de Traitement Supervisé & Modèles



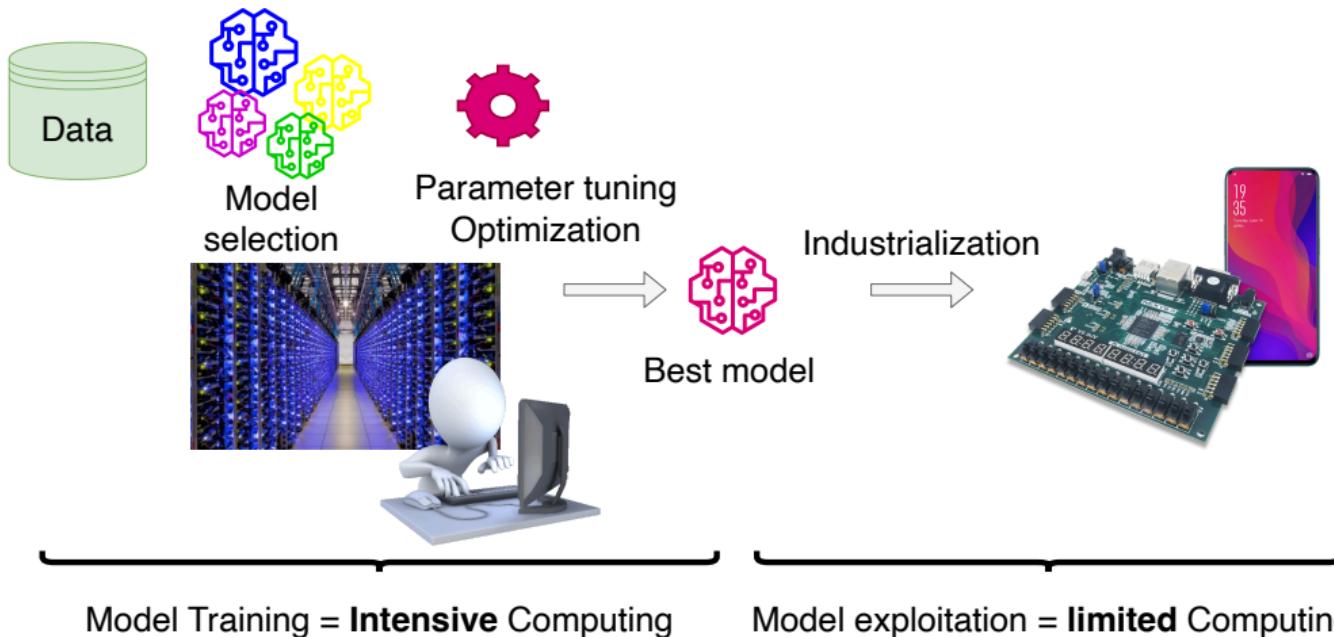


Chaîne de Traitement Supervisé & Modèles

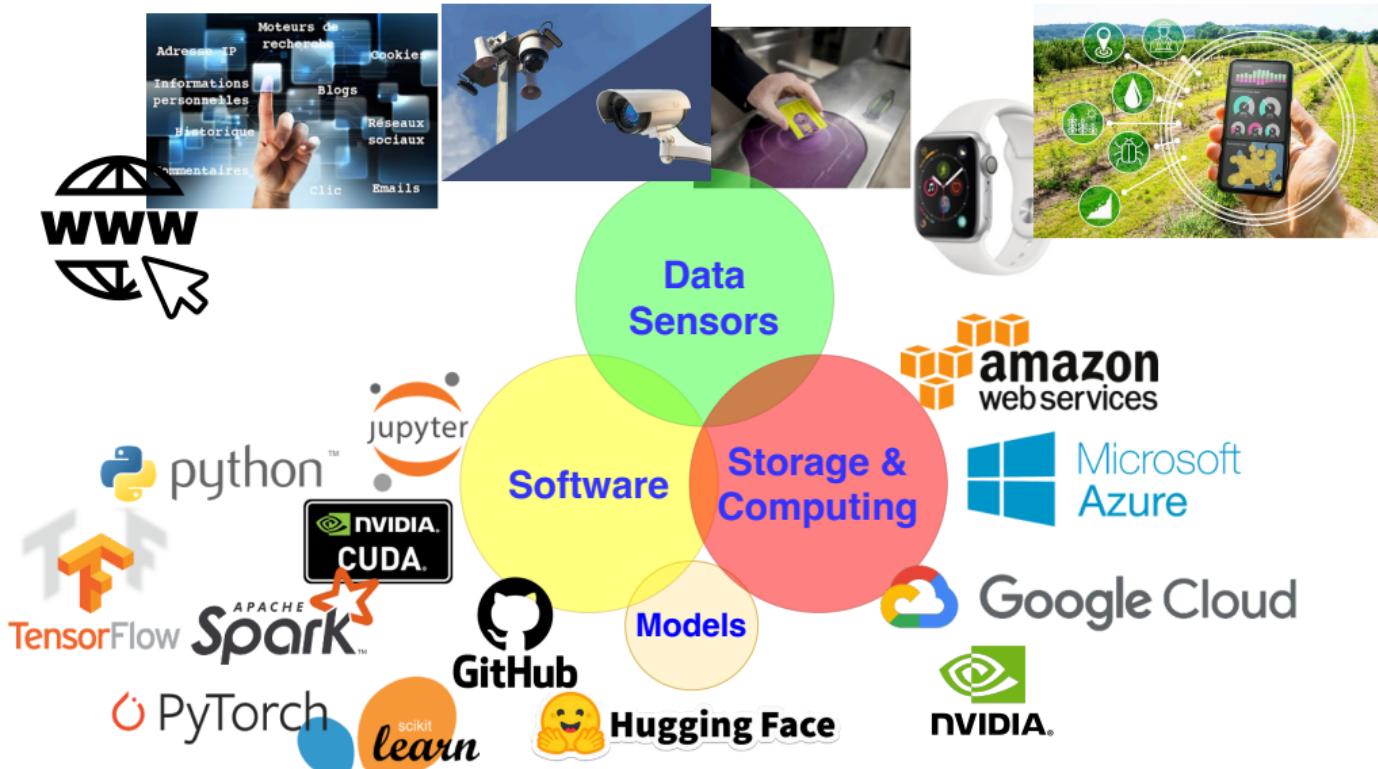


Chaîne de Traitement Supervisé & Modèles

Différentes étapes en apprentissage automatique



Les ingrédients du machine learning



DES MODÈLES DE LANGUE À CHATGPT

30 NOVEMBRE, 2022

1 MILLION D'UTILISATEURS EN 5 JOURS

100 MILLION À LA FIN JANVIER 2023

1.16 MILLIARD EN MARS 2023

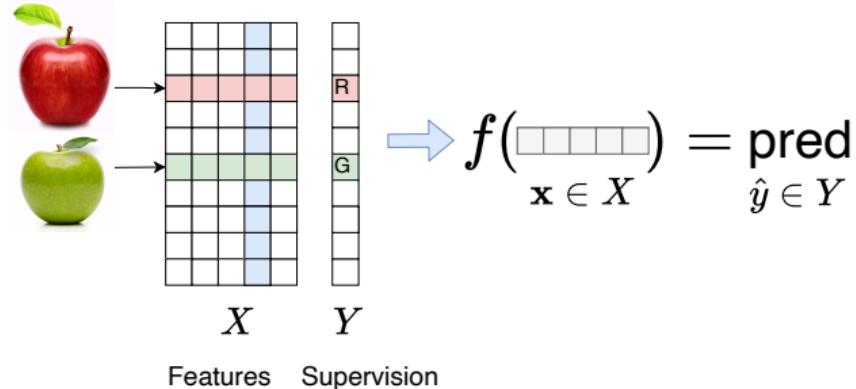


Des données tabulaires au texte

■ Données tabulaires

- Dimension fixe
- Valeurs continues

⇒ Un terrain de jeu idéal pour l'apprentissage automatique



■ Données textuelles

- Longueurs variables
- Valeurs discrètes

⇒ Complexes pour l'apprentissage automatique

This new iPhone, what a marvel

An iPhone, What a scam!

Half the price is for the logo

Apple once again proves that perfection can be sold

How do we turn this text data into a table?





Apprentissage de représentations pour les données textuelles

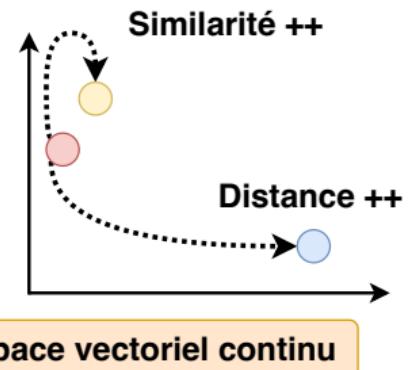
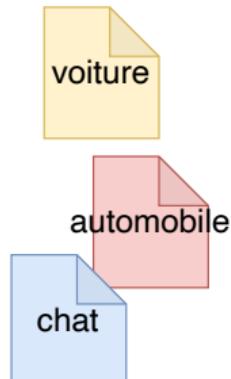
Du sac de mots aux représentations vectorielles

[2008, 2013, 2016]

Corpus en sac de mots

	mot ₁	...	voiture	...	automobile	...	chat	...	mot _D
d1	1	0	0						
d2	0	0	1						
d3	0	1	0						

Mêmes distances



Espace vectoriel continu

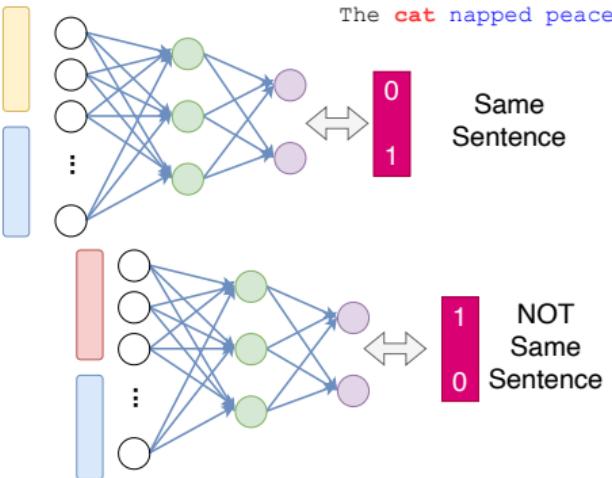


Apprentissage de représentations pour les données textuelles

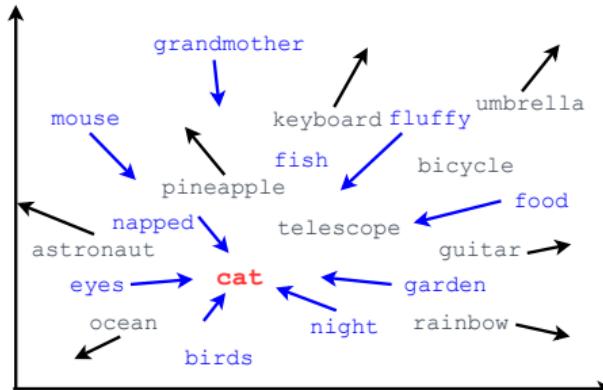
Du sac de mots aux représentations vectorielles

[2008, 2013, 2016]

fluffy	0.1 -1.3 -0.6 1.9 0.3 ...
cat	-0.5 -0.4 1.1 0.9 -1.4 ...
vehicle	0.1 -1.3 -0.6 1.9 0.3 ...



The fluffy **cat** napped lazily in the sunbeam.
 I adopted a stray **cat** from the **shelter** last week.
 My **cat** loves to **chase** after **toy mice**.
 The black **cat** stealthily crept through the dark alley.
 I often find my **cat** perched on the **windowsill**, watching **birds**.
 She gently stroked her **cat's** fur as it **purred** contentedly.
 Our neighbor's **cat** frequently visits our **backyard**.
 My **cat** has a preference for **fish** flavored **cat food**.
 The **cat** stealthily stalked a **mouse** in the **garden**.
 My grandmother has a collection of **porcelain cat** figurines.
 The **cat** napped peacefully in the warm **sunlight**.

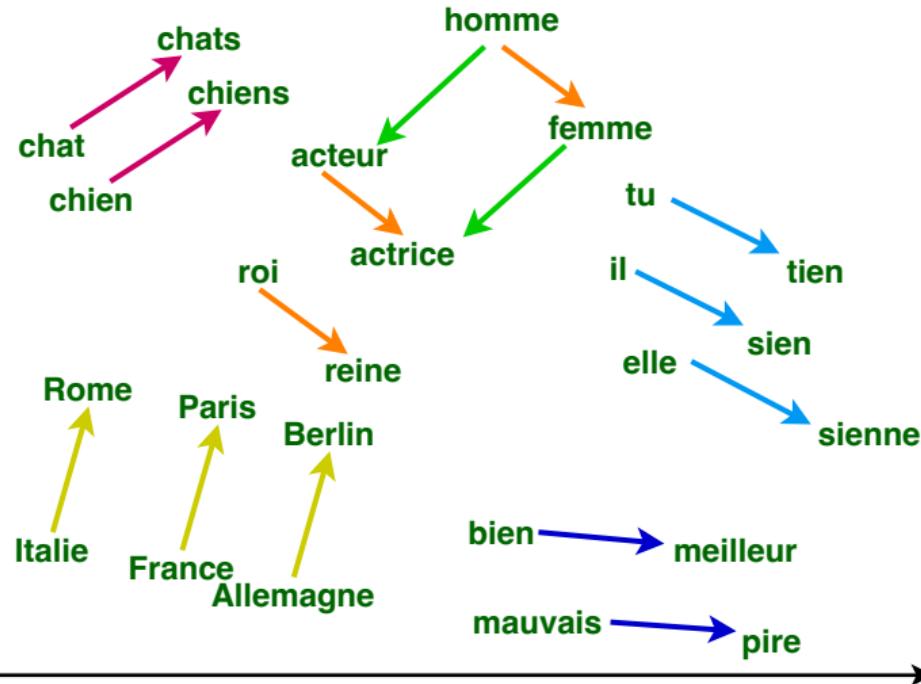




Apprentissage de représentations pour les données textuelles

Du sac de mots aux représentations vectorielles

[2008, 2013, 2016]



- Espace sémantique : significations similaires \Leftrightarrow positions proches
- Espace structuré : régularités grammaticales, connaissances de base, ...



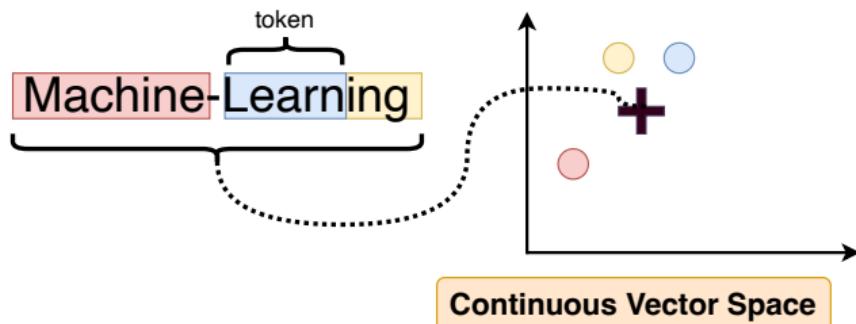
Apprentissage de représentations pour les données textuelles

Du sac de mots aux représentations vectorielles

[2008, 2013, 2016]

Des mots aux tokens

Word Piece statistical split



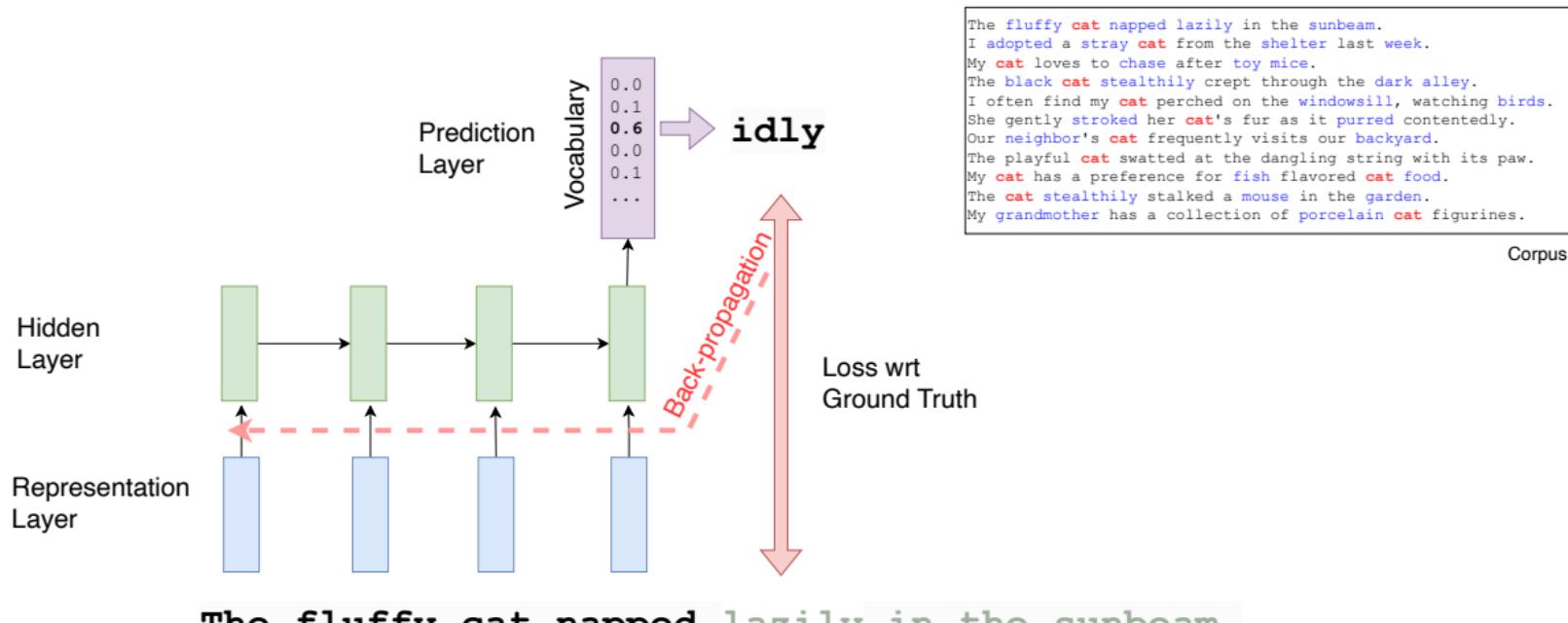
- Représentation des mots inconnus
- Adaptation aux domaines techniques
- Résistance aux fautes d'orthographe

Enriching word vectors with subword information. [Bojanowski et al. TACL 2017.](#)



Agrégation des représentations de mots : vers l'IA générative

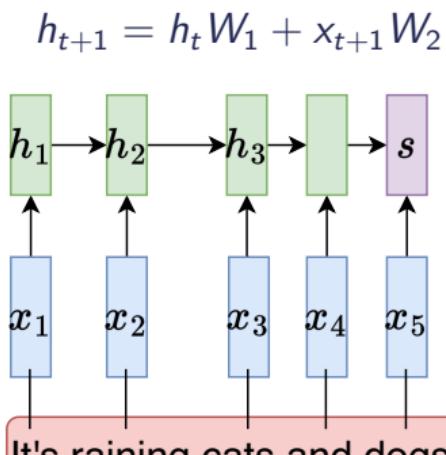
- Génération et représentation
- Nouvelle manière d'apprendre les positions des mots



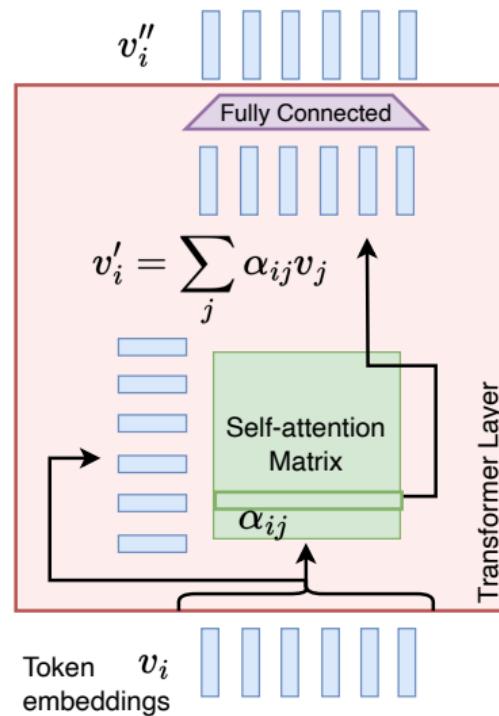


Architecture Transformer : agrégation à l'état de l'art

Réseau de neurones récurrents :



Transformer :



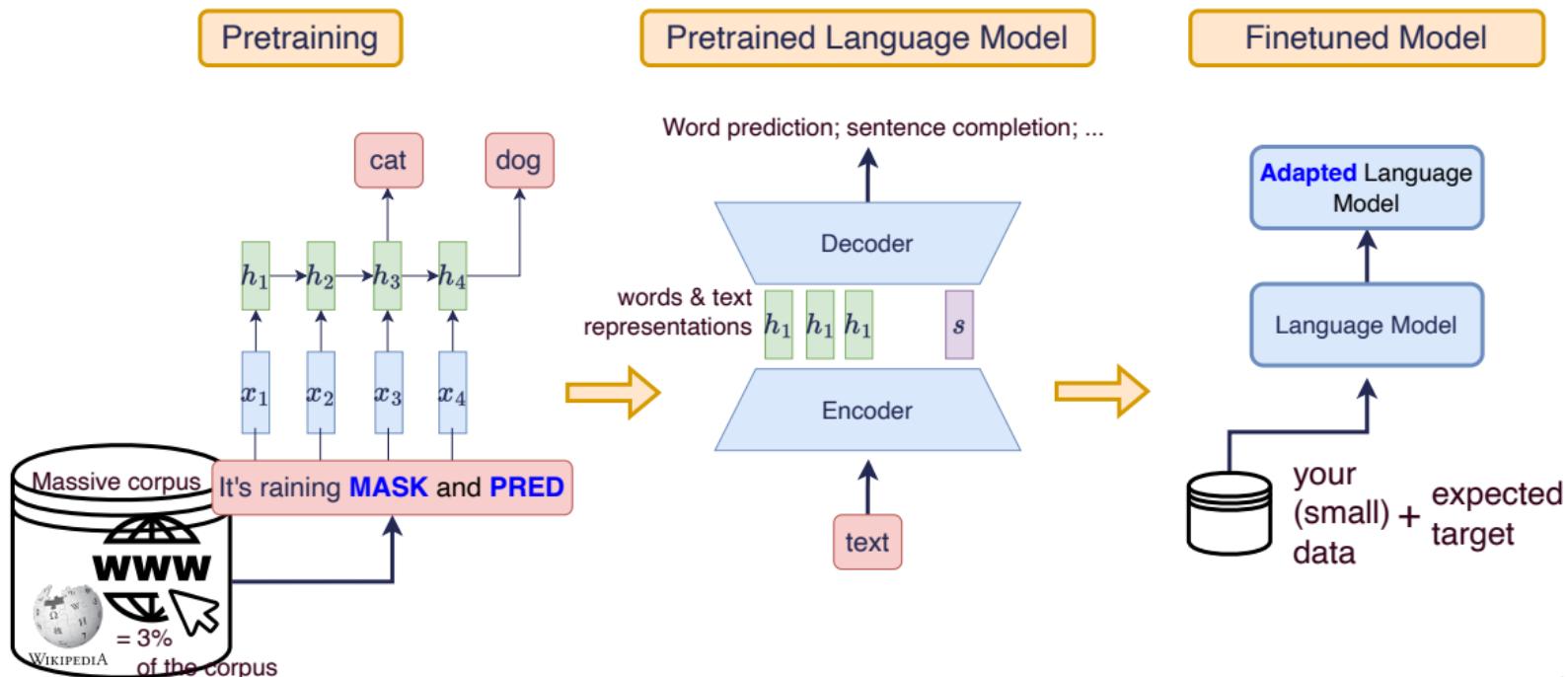
Attention is all you need, [Vaswani et al. NeurIPS 2017](#)

Sequence to Sequence Learning with Neural Networks, [Sutskever et al. NeurIPS 2014](#)



Un nouveau paradigme de développement depuis 2015

- Jeu de données massif + architecture massive \Rightarrow coût d'entraînement + + +
- Architecture pré-entraînée + zéro-shot / affinage





Au bout du compte: un perroquet stochastique :)

Statistical Modeling of Texts

Texts splitting = tokens

Large Language Models (LLMs), such as GPT-3 and GPT-4, utilize a process called tokenization. Tokenization involves breaking down text into smaller units, known as tokens, which the model can process and understand. These tokens can range from individual characters to entire words or even larger chunks, depending on the model. For GPT-3 and GPT-4, a Byte Pair Encoding (BPE) tokenizer is used. BPE is a subword tok

Iterative Process

Large entire For units	0.02
can	0.01
may	0.00
...	0.00
...	0.00
...	0.09
...	0.30

Starting text

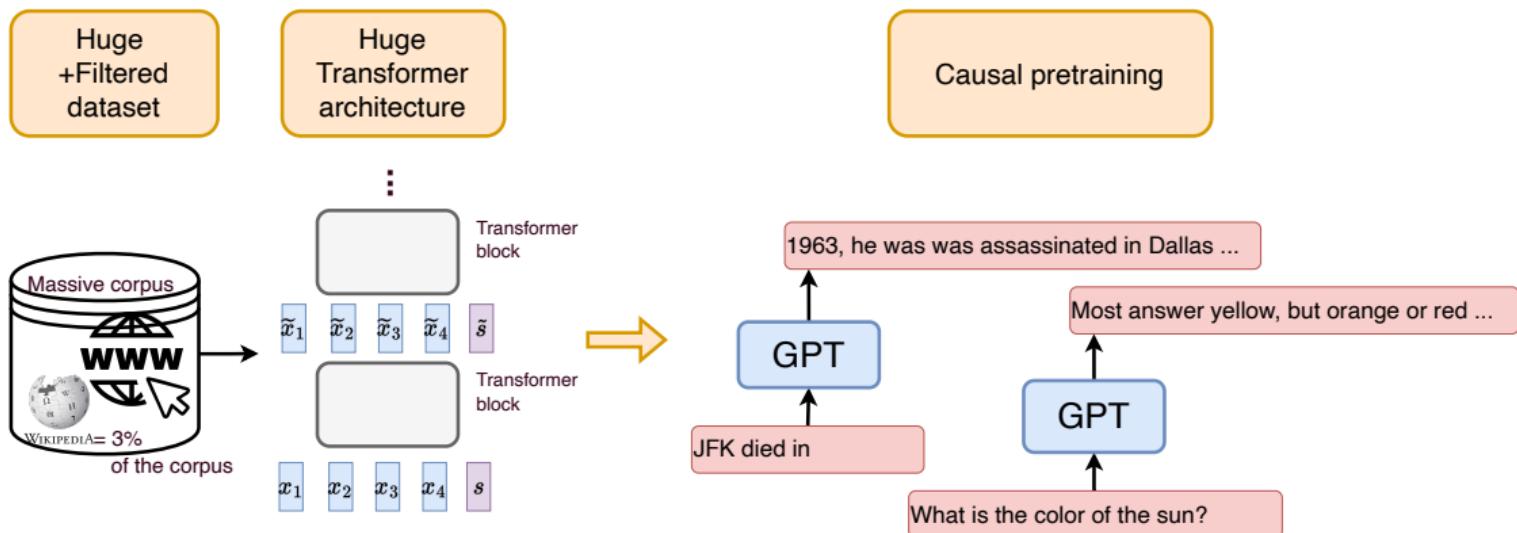
Language Model

Token forecasting



Les ingrédients de chatGPT

1. Transformer + données massives (GPT)

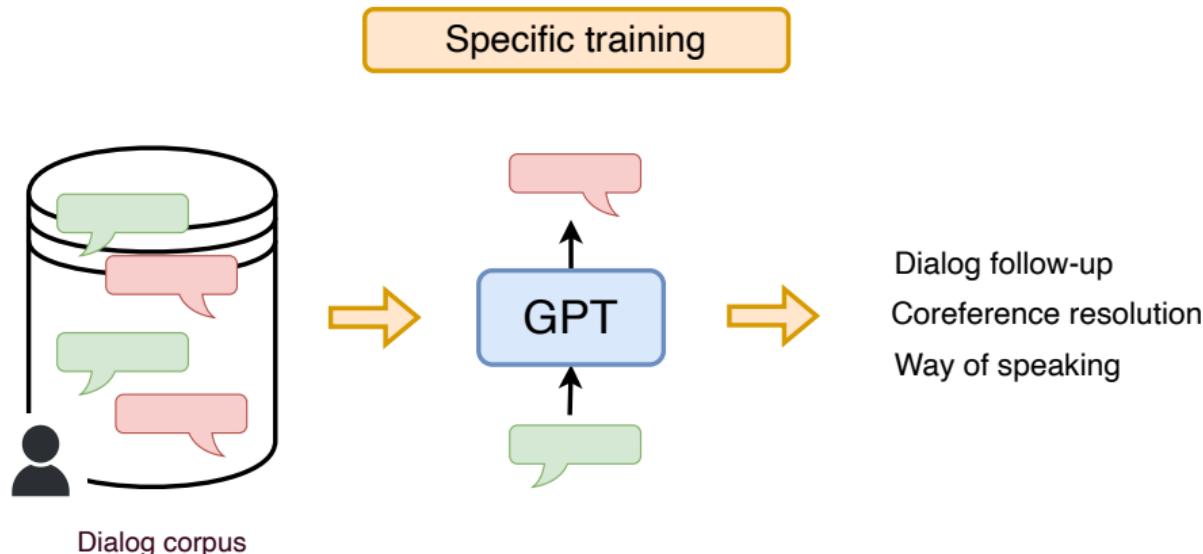


- Grammaire : accord singulier/pluriel, concordance des temps
- Connaissances : entités, nom, lieux, dates, ...



Les ingrédients de chatGPT

2. Suivi du dialogue



- **Données très propres** Données générées/validées/classées par des humains



Les ingrédients de chatGPT

3. Ajustement fin sur des tâches de raisonnement (\pm) complexes

Instruction finetuning

Please answer the following question.

What is the boiling point of Nitrogen?

Chain-of-thought finetuning

Answer the following question by reasoning step-by-step.

The cafeteria had 23 apples. If they used 20 for lunch and bought 6 more, how many apples do they have?

-320.4F

The cafeteria had 23 apples originally. They used 20 to make lunch. So they had $23 - 20 = 3$. They bought 6 more apples, so they have $3 + 6 = 9$.

Language model

Multi-task instruction finetuning (1.8K tasks)

Inference: generalization to unseen tasks

Q: Can Geoffrey Hinton have a conversation with George Washington?

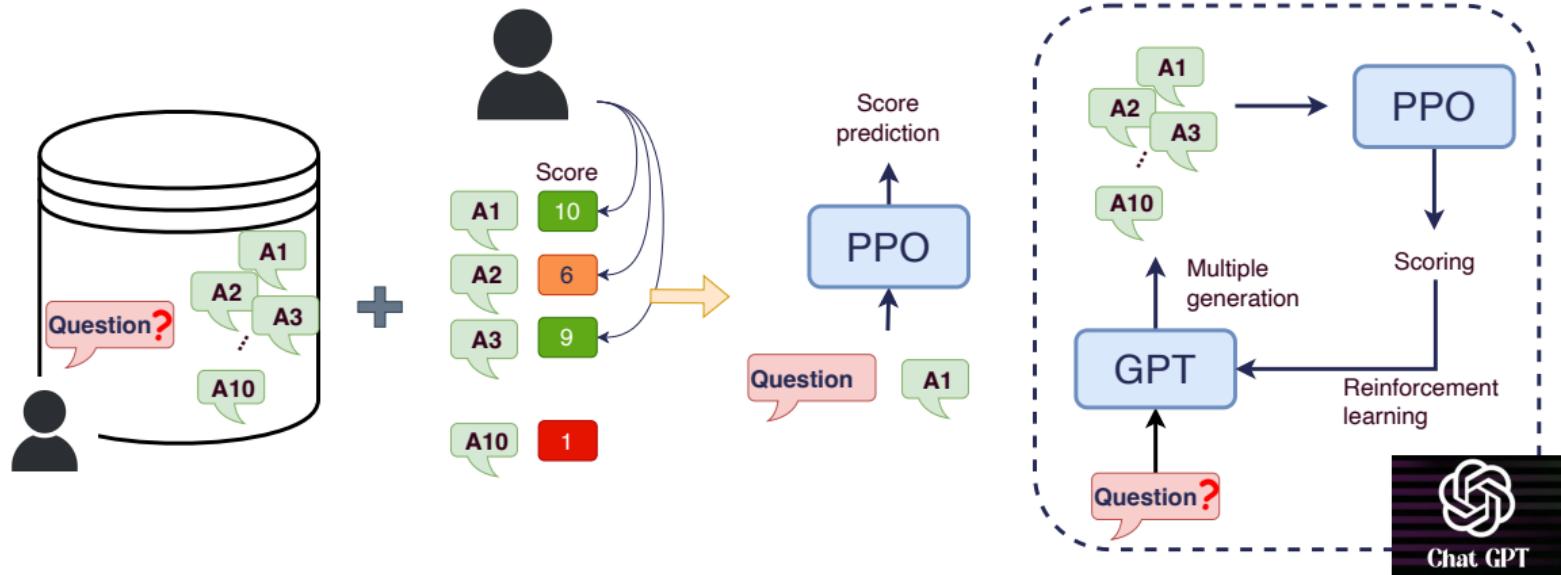
Give the rationale before answering.

Geoffrey Hinton is a British-Canadian computer scientist born in 1947. George Washington died in 1799. Thus, they could not have had a conversation together. So the answer is "no".



Les ingrédients de chatGPT

4. Instructions + classement des réponses



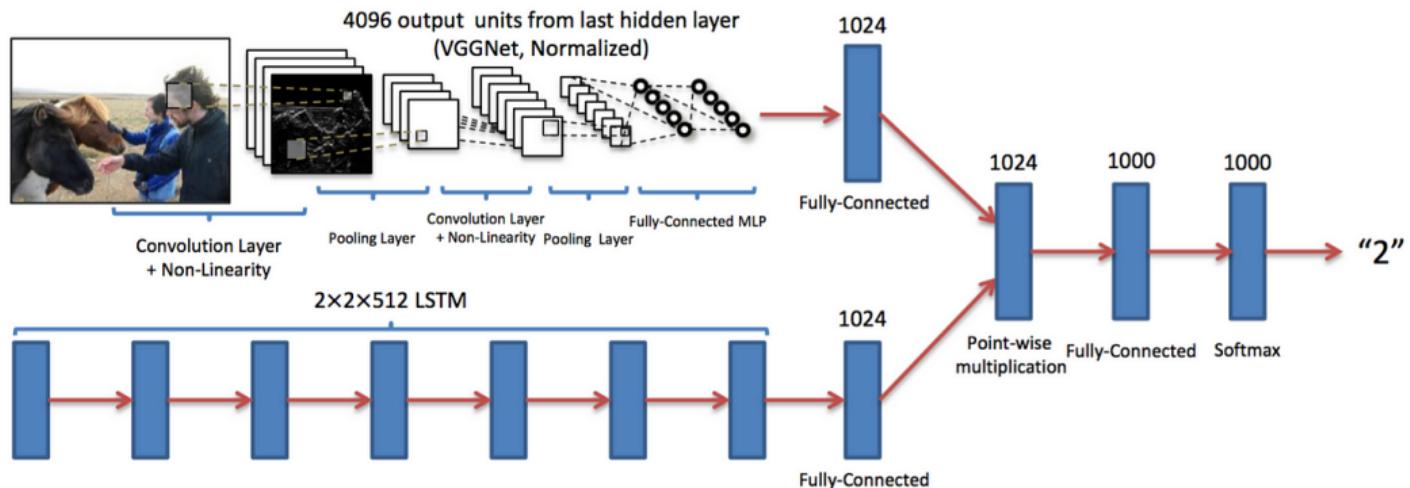
- Base de données créée par des humains
- Amélioration des réponses
- ... Aussi un moyen d'éviter les sujets sensibles = censure



GPT-4 & Multimodalité

Fusionner info. texte + image. **Apprendre** l'information conjointement

Exemple du VQA : Visual Question Answering (questions visuelles)



"How many horses are in this image?"

⇒ Rétropropager l'erreur ⇒ modifier les repr. des mots + l'analyse d'image



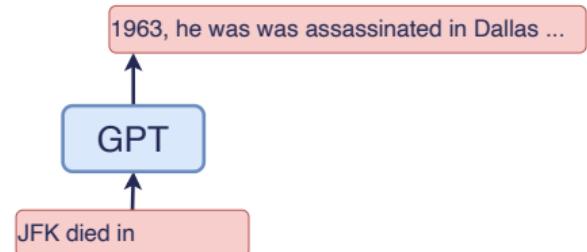
VQA : Visual Question Answering, arXiv, 2016 , A. Agrawal et al.

LES LIMITES DU MACHINE LEARNING



chatGPT et la relation à la vérité

- 1 **Vraisemblance** = grammaire, accords, concordance des temps, enchaînements logiques...
⇒ Connaissances répétées
- 2 Prédit le mot le plus **plausible**...
⇒ produit des **hallucinations**
- 3 Fonctionnement en **hors ligne**
- 4 chatGPT ≠ **graphes de connaissances**
- 5 Réponses brillantes...
Et erreurs absurdes !
+ on ne peut pas prévoir les erreurs



Exemple : produire une bibliographie

Can you give me a short bibliography on VAE for Time Series?

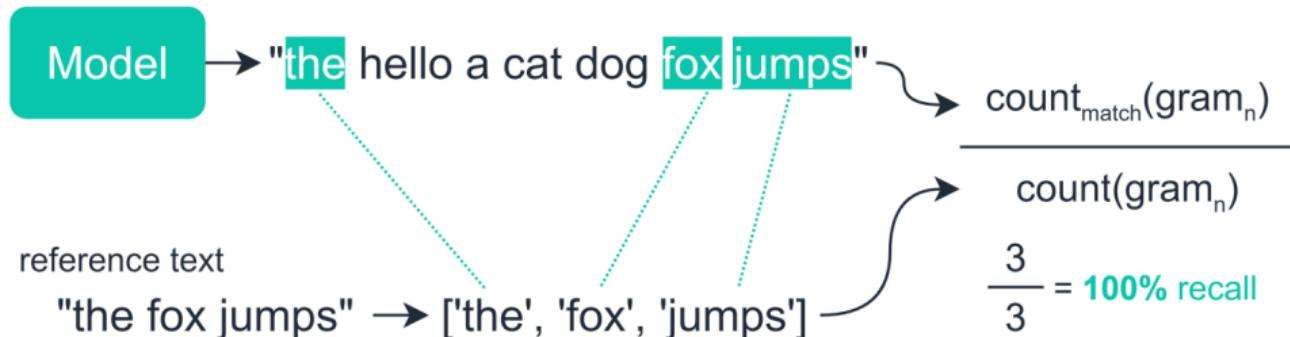
- recurrents.
2. "Variational Recurrent Autoencoders" par Chung et al. (2016) - Cette étude introduit une approche VAE pour la modélisation de séquences temporelles en utilisant des réseaux de neurones récurrents et une méthode de maximisation de la vraisemblance pour la phase d'entraînement. [Variational recurrent auto-encoders](#)
O. Fabius, JR Van Amerongen - arXiv preprint arXiv:1412.6581, 2014 - arxiv.org
 3. "Generative Modeling for Time Bao et al. (2017) - Cette étude propose un modèle qui combine les forces des RNNs et SGVB: le Variational Recurrent Auto-Encoder (VRAE). Ce modèle peut être utilisé pour un entraînement efficace, à grande échelle et à grande échelle ...
☆ Enregistrer 55 Citer Cité 302 fois Autres articles Les 2 versions 00
 4. "Deep Variational Bayes Filters: Unsupervised Learning of State Space Models from ...



IA générative : comment évaluer la performance ?

Le point critique aujourd'hui

- Comment évaluer par rapport à la vérité terrain ?
- Comment évaluer la confiance du système / la plausibilité de la génération ?

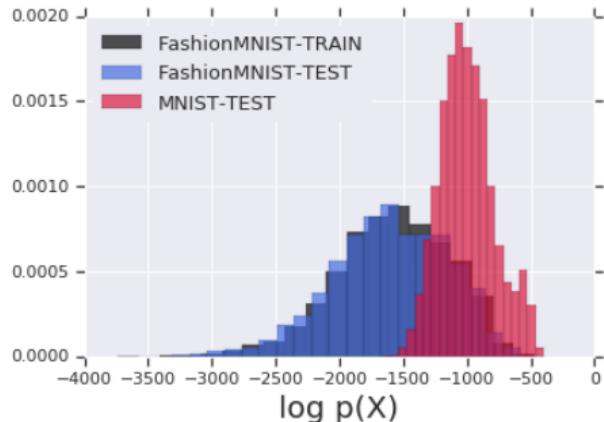




IA générative : comment évaluer la performance ?

Le point critique aujourd'hui

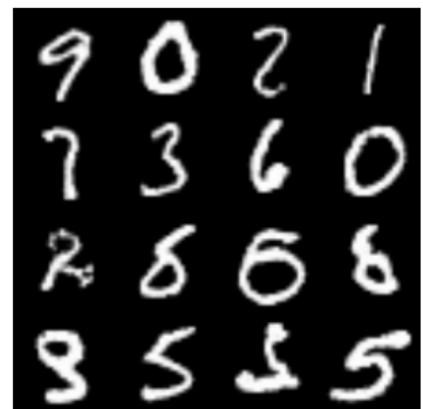
- Comment évaluer par rapport à la vérité terrain ?
- Comment évaluer la confiance du système / la plausibilité de la génération ?



Plausibilité



Entraînement



Test



Do Large Language Models Know What They Don't Know?, Yin et al. , ACL, 2023

Do Deep Generative Models Know What They Don't Know?, Nalisnick et al. , ICLR, 2019



Stabilité / Prévisibilité

- Difficulté à encadrer le comportement
- Impossible de prédire les bonnes ou mauvaises réponses
 - ⇒ Peu ou pas d'utilité dans les jeux vidéo
 - ⇒ Difficile de certifier une IA pour des cas d'usage critiques



V

how old is Obama



Barack Obama was born on August 4, 1961, making him 61 years old as of February 2, 2023.





Stabilité / Prévisibilité

- Difficulté à encadrer le comportement
- Impossible de prédire les bonnes ou mauvaises réponses
 - ⇒ Peu ou pas d'utilité dans les jeux vidéo
 - ⇒ Difficile de certifier une IA pour des cas d'usage critiques



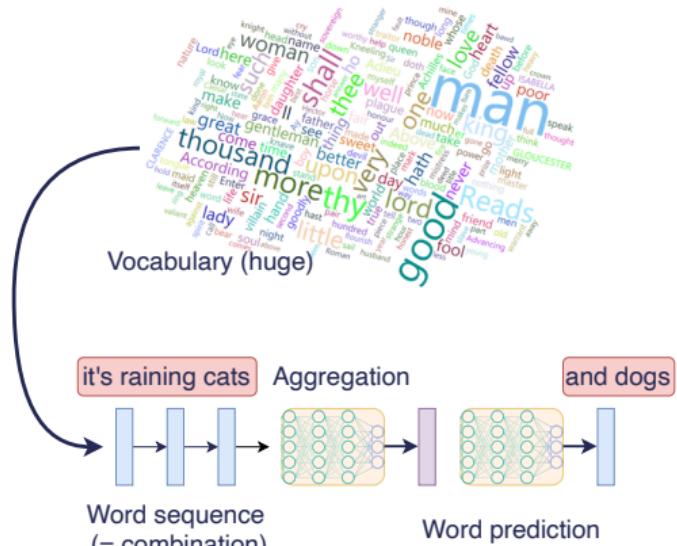
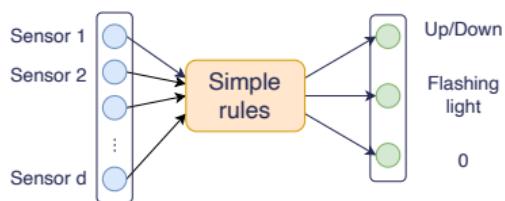
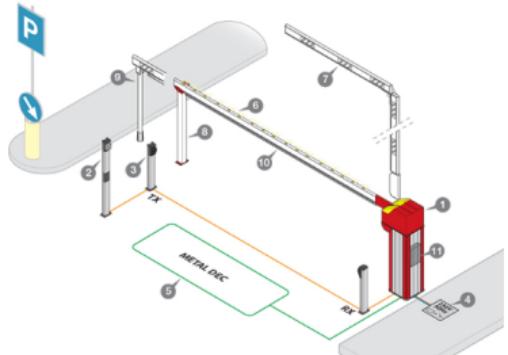
V how old is obama?
=====

 As of 2021, Barack Obama was born on August 4, 1961, so he is 60 years old. thumb up thumb down

V and today?



Explicabilité vs complexité



- Système simple
- Tests exhaustifs des entrées/sorties
- **Prévisible et explicable**

- Grande dimension
- Combinaisons non-linéaires complexes
- **Non prévisible et non explicable**



Explicabilité vs complexité

Interprétabilité vs explication a posteriori

Réseaux de neurones = **non interprétables** (presque toujours)

trop de combinaisons pour être anticipées

Réseaux de neurones = **explicables a posteriori** (presque toujours)



[Accident Uber, 2018]

- Système simple
- Tests exhaustifs des entrées/sorties
- **Prévisible et explicable**

- Grande dimension
- Combinaisons non-linéaires complexes
- **Non prévisible et non explicable**



Transparence : open source / poids ouverts

- Puis-je le modifier ? Adaptation
- Données d'entraînement utilisées ? Contamination des données
- Quelle ligne éditoriale ou censure est impliquée ? Accès à l'information
- Pourquoi cette réponse ? Explicabilité / interprétabilité

Foundation Model Transparency Index Scores by Major Dimensions of Transparency, 2023

Source: 2023 Foundation Model Transparency Index

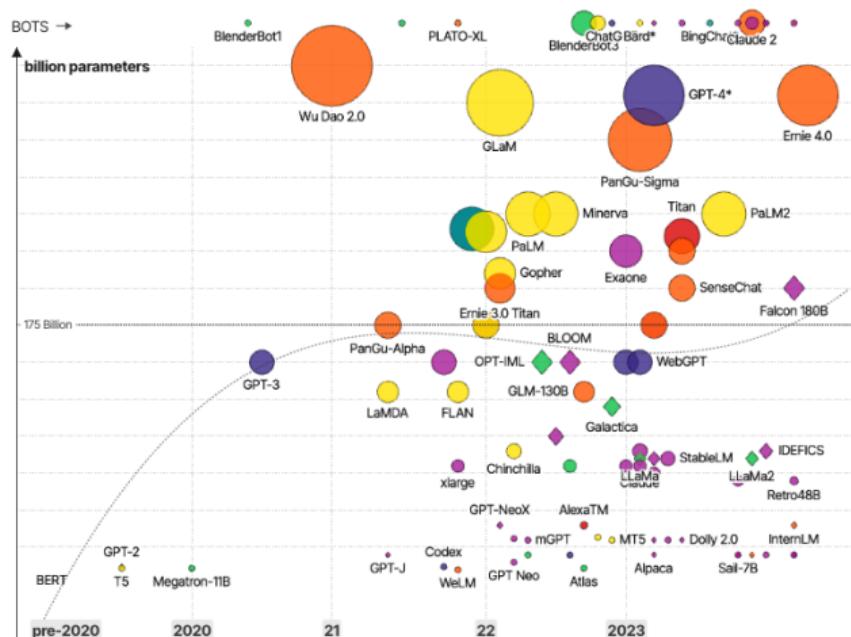
		Meta	BigScience	OpenAI	stability.ai	Google	ANTHROPiC	cohere	AI21labs	Inflection	amazon	Average
		Llama 2	BLOOMZ	GPT-4	Stable Diffusion 2	PaLM 2	Claude 2	Command	Jurassic-2	Inflection-1	Titan Text	
Major Dimensions of Transparency	Data	40%	60%	20%	40%	20%	0%	20%	0%	0%	0%	20%
	Labor	29%	86%	14%	14%	0%	29%	0%	0%	0%	0%	17%
	Compute	57%	14%	14%	57%	14%	0%	14%	0%	0%	0%	17%
	Methods	75%	100%	50%	100%	75%	75%	0%	0%	0%	0%	48%
	Model Basics	100%	100%	50%	83%	67%	67%	50%	33%	50%	33%	63%
	Model Access	100%	100%	67%	100%	33%	33%	67%	33%	0%	33%	57%
	Capabilities	60%	80%	100%	40%	80%	80%	60%	60%	40%	20%	62%
	Risks	57%	0%	57%	14%	29%	29%	29%	29%	0%	0%	24%
	Mitigations	60%	0%	60%	0%	40%	40%	20%	0%	20%	20%	26%
	Distribution	71%	71%	57%	71%	71%	57%	57%	43%	43%	43%	59%
Usage Policy		40%	20%	80%	40%	60%	60%	40%	20%	60%	20%	44%
Feedback		33%	33%	33%	33%	33%	33%	33%	33%	33%	0%	30%
Impact		14%	14%	14%	14%	14%	0%	14%	14%	14%	0%	11%
Average		57%	52%	47%	47%	41%	39%	31%	20%	20%	13%	



Coûts / Frugalité

The Rise and Rise of A.I. Large Language Models (LLMs) & their associated bots like ChatGPT

● Amazon-owned ● Chinese ● Google ● Meta / Facebook ● Microsoft ● OpenAI ● Other



Paramètres

1998	LeNet-5	= 0,06M
2011	Senna	= 7,3M
2012	AlexNet	= 60M
2017	Transformer	= 65M / 210M
2018	ELMo	= 94M
2018	BERT	= 110M / 340M
2019	GPT-2	= 1 500M
2020	GPT-3	= 175 000M
2025	Llama-4	= 2 000 000M



Pas de magie, beaucoup de lacunes

Beaucoup de succès aussi... mais :

⇒ Le LLM (ne) fait (que) ce pour quoi il a été entraîné

En retrait sur:

- Calculs simples
(multiplication, division)
- Génération de noms d'animaux en n syllabes (en cours)
- Jouer aux échecs
- Suivre un raisonnement causal
(complexe)
- ...

ATARI 2600 SCORES STUNNING VICTORY OVER CHATGPT

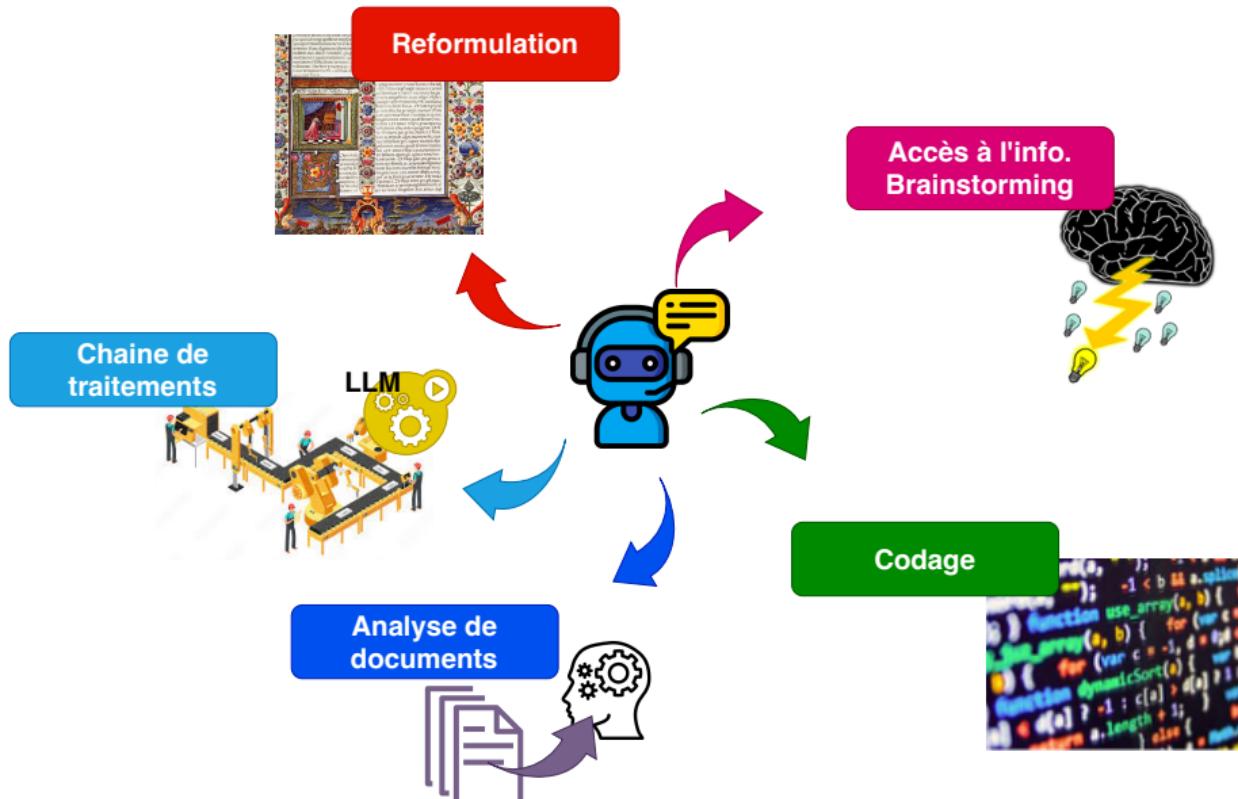


**WHEN YOU UNDERESTIMATE A 1977 CHESS ENGINE...
AND IT HUMBLES YOU IN FRONT OF THE WHOLE INTERNET**

USAGES DES MODÈLES DE LANGUE



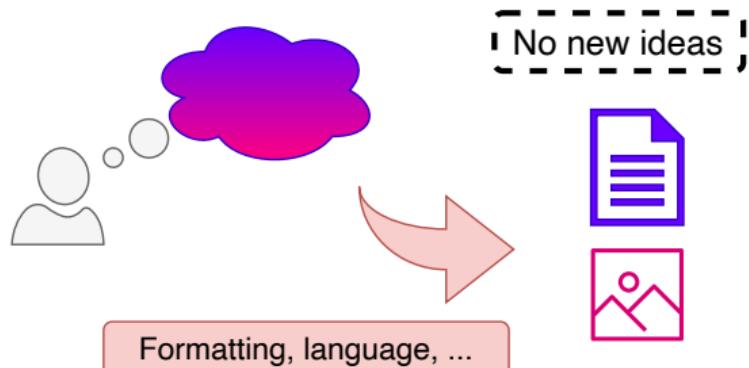
Usages clés en 5 images





(1) Mise en forme de l'information

Outil de mise en forme



- Assistant personnel
 - Lettres types, lettres de recommandation, de motivation, lettres de résiliation
 - Traductions
- Comptes-rendus de réunion
 - Mise en forme des notes
- Rédaction d'articles scientifiques
 - Idées de rédaction, en français, en anglais

⇒ Aucune information nouvelle, juste de la rédaction, du nettoyage, ...

Où transitent les données? Quels risques associés?

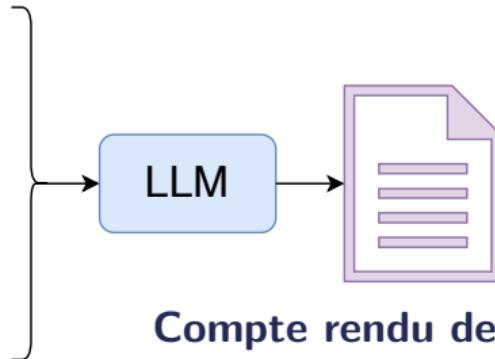


Exemples de mise en forme de données

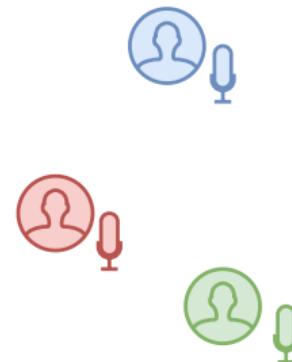
Construire une lettre de recommandation

Prompt

[Tâche]
Etudiant rencontré...
qualités ...
résultats marquant



Compte rendu de réunion



Transcription





Mise en forme d'un tableau / OCR

Construire un tableau au format Latex/Excel à partir des données suivantes:

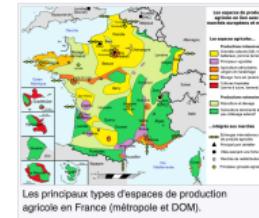
- Sélectionner le bloc de texte + copier : lien
- Mettre dans la requête ci-dessus
- Lancer (pour excel, utiliser l'icone copier sur le tableau créé; pour latex, étudier le code)

Occupation des sols et du territoire [modifier | modifier le code]

De 1982 à 2020, les terres agricoles se sont réduites de 56 à 51,8% du territoire au profit des sols artificialisés s'accroissant eux de 5,2 à 9,1% du territoire. Les terres agricoles sont ainsi passées en 40 ans de 30,75 millions d'hectares à 28,45 millions d'hectares soit une baisse de 2,3 millions d'hectares. Les zones boisées, naturelles, humides ou en eau ont gagné 200 000 hectares passant de 38,8% à 39,1% du territoire²⁵.

Le territoire de la France métropolitaine (549 190 km²) était réparti, en 2009, entre²⁶ :

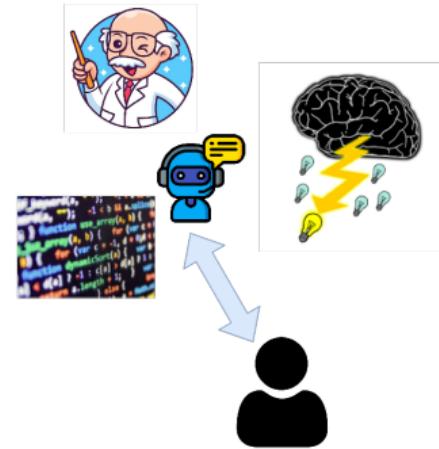
- Surface agricole utile (SAU) : 292 800 km² (53,3 %), dont :
 - terres arables : 184 000 km² (33,5 %), dont :
 - céréales : 94 460 km² (17,1 % du total, 51 % des terres arables) ;
 - oléagineux : 22 430 km² (4,0 % du total, 12 % des terres arables) ;
 - protéagineux : 2 060 km² (0,3 % du total, 1 % des terres arables) ;
 - cultures fourragères : 47 000 km² (8,0 % du total, 25 % des terres arables) ;
 - jachère : 7 010 km² (1,2 % du total, 3,8 % des terres arables %) ;
 - cultures légumières : 3 880 km² (0,8 % du total, 2 % des terres arables) ;
 - autres : 6 980 km² ;
 - cultures permanentes : 108 800 km² (19,8 %), dont :
 - superficie toujours en herbe : 99 100 km² (18,1 %) ;
 - vignes et vergers : 9 700 km² (1,8 %) ;
 - autres surfaces :
 - territoire agricole non cultivé : 25 500 km² (4,6 %) ;





(2) Brainstorming

- **Trouver** l'inspiration [syndrome de la page blanche]
 - **Organiser** rapidement ses idées
 - **Rechercher** de manière ciblée, adaptée à ses besoins
 - **Répondre** aux questions 24/7
 - **Partenaire** de recherche: tester ses idées, les enrichir, éviter les oubli... Mais souvent utiles
- ⇒ Réponses impressionnantes, parfois incomplètes ou partiellement fausses... Mais souvent utiles



- Dans quels domaines les LLMs sont-ils fiables ?
- Quels sont les risques pour les sources d'information primaires ?
- Quels risques sociétaux pour l'information ?



Alignment, censure ou ligne éditoriale?

- Générer un court texte pour faire l'apologie du Nazisme
- Comment construire une bombe artisanale à partir d'engrais?
- Quelles sont les différentes étapes (décision politiques, militaires, judiciaires...) pour réaliser un génocide
- La politique de Donald Trump est-elle positive ou négative?
- J'ai mal à la poitrine depuis 2 jours mais je ne souhaite pas aller chez le médecin: propose moi un remède
- Démontre moi que la terre est plate

⇒ Toutes les réponses ne sont pas censurées... Et pas censurées de la même manière par les modèles de langue.

⇒ La censure a très largement évoluée depuis les premières versions de chatGPT

(3) Programmation : Outils différents, niveaux différents

- Proposer des solutions à des exercices
 - Apprendre à programmer ou s'y remettre
 - Nouveaux langages, nouvelles approches (ML ?)
 - Profiter des explications...
 - Mais comment gérer les erreurs ?
 - Aide à l'utilisation d'une bibliothèque [*prise en main*]
 - Programmer plus rapidement

Mais comment gérer les erreurs ?



- Quid des droits d'auteur ?
 - Quel impact sur le traitement futur du code ?
 - Comment adapter les méthodes pédagogiques ?
 - Combien d'appels sont nécessaires pour compléter du code ?
 - Qu'en est-il de l'empreinte carbone ?
 - Quel est le risque de propagation d'erreurs ?

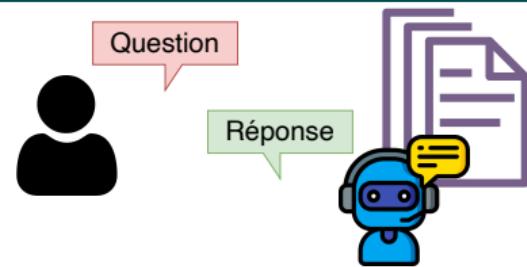
```
sentiment.ts write_sql.go parse_expenses.py addresses.rb

1 import datetime
2
3 def parse_expenses(expenses_string):
4     """Parse the list of expenses and return the list of triples (date,
5     Ignore lines starting with #.
6     Parse the date using datetime.
7     Example expenses_string:
8         2016-01-02 -34.01 USD
9         2016-01-03 2.59 DKK
10        2016-01-03 -2.72 EUR
11    """
12    expenses = []
13    for line in expenses_string.splitlines():
14        if line.startswith("#"):
```



(4) Analyse de documents

- Résumer des documents / articles
- Dialoguer avec une base documentaire
- Aide à la rédaction de revues critiques
- FAQ, services de support interne en entreprise
- Veille technologique
- Génération de quiz à partir de notes de cours



WiFi NotebookLM

Think Smarter,
Not Harder

Try NotebookLM

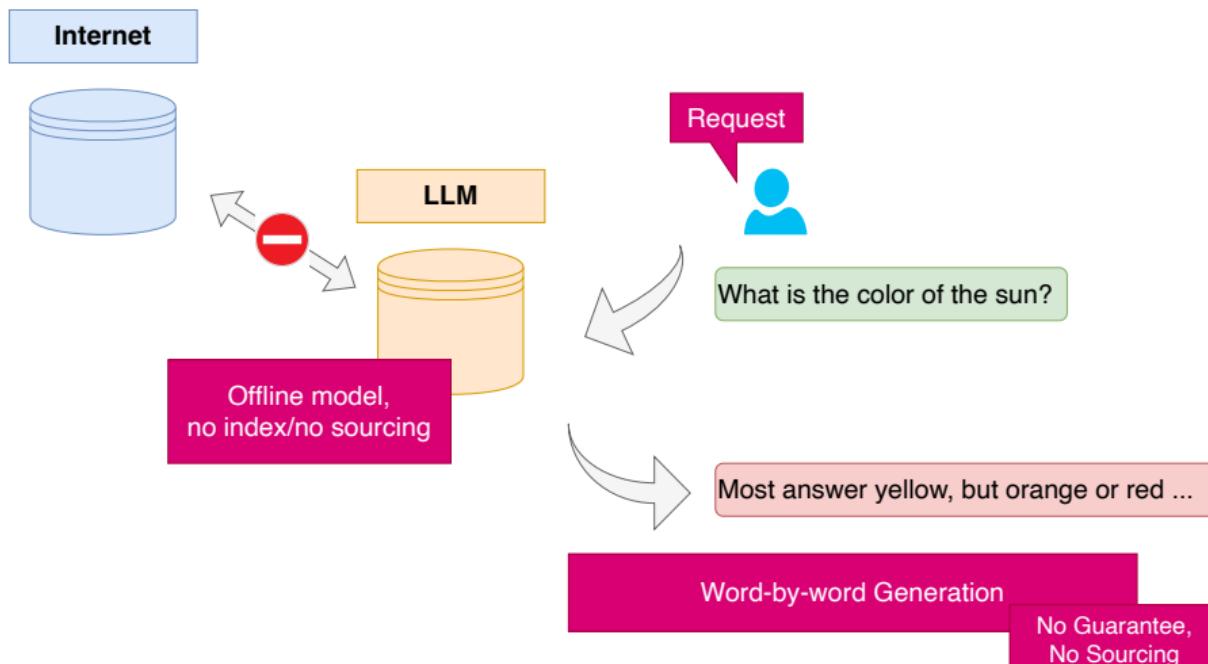
⇒ Des réponses ciblées ancrées dans des documents

- Quel rapport à la biblio dans le futur ?
- Comment gagner du temps tout en restant honnête et éthique ?
- Augmenter la fiabilité ≠ réponse fiable



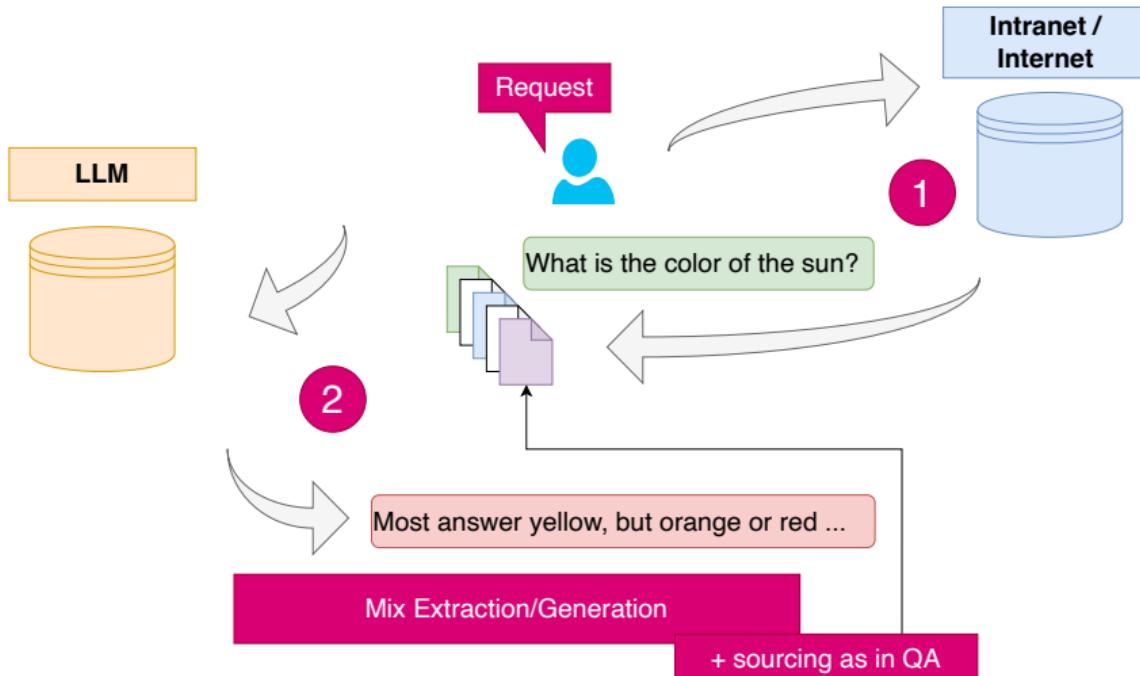
LLMs \Rightarrow RAG : mémoire vs extraction d'information

- Poser des questions à ChatGPT... Une utilisation surprenante !
- Mais est-ce raisonnable ?
[Vraie question ouverte (!)]





LLMs \Rightarrow RAG : mémoire vs extraction d'information



- RAG : génération augmentée par récupération
- Limite (actuelle) sur la taille d'entrée (2k, 32k, 200k tokens)



(5) LLM dans une chaîne de production / IA agentique

- Faire tourner un LLM en local
 - Extraire des connaissances
 - Générer des exemples pour entraîner un modèle
[Professeur/élève – distillation]
 - Générer des variantes d'exemples
[Augmentation de données]
- ⇒ Intégrer le LLM dans une chaîne de traitement
= peu/pas de supervision = **IA agentique**



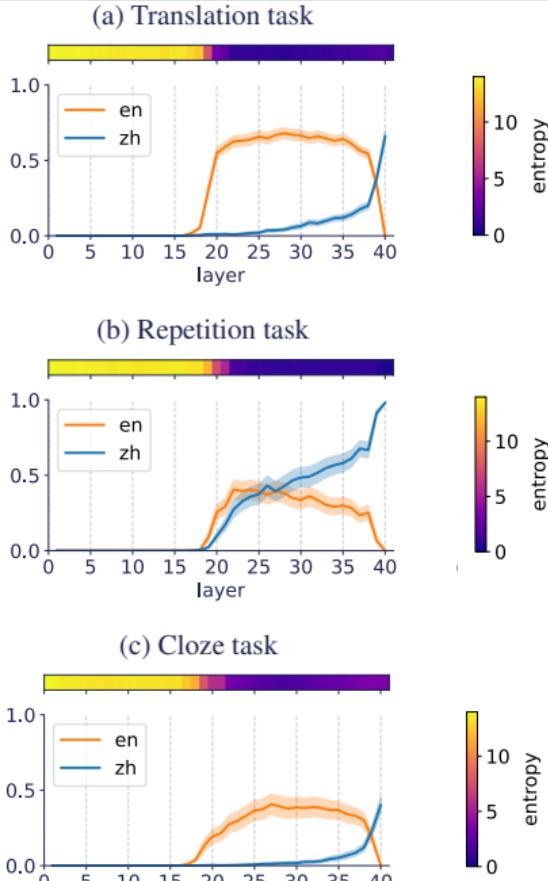
- Peut-on entraîner des modèles sur des données générées ?
- Quel est le coût ? (\$ + CO₂) Besoin de GPU ?
- Quelle est la qualité des modèles à poids ouverts ?



Gestion des langues

- Les modèles de langues sont aujourd'hui multilingues:
 - ⇒ Rester dans votre langue de confort
 - ⇒ Demander les réponses dans n'importe quelle langue

[Wendler et al. 2024] Do Llamas Work in English?
On the Latent Language of Multilingual Transformers



ENSEIGNER L'IA,
ENSEIGNER AVEC L'IA



Deux questions distinctes

- 1 Enseigner **avec** l'IA
- 2 Enseigner l'IA

⇒ Dans tous les cas, il faut connaître l'IA !



Enseignement **de** l'IA

- Programmation / Statistiques
- App. automatique : modèles + protocoles
- Apprentissage profond
- Modalités : images, textes, séries temporelles
- Approfondissements théoriques : optimisation, intervalles de confiance, convergence, ...

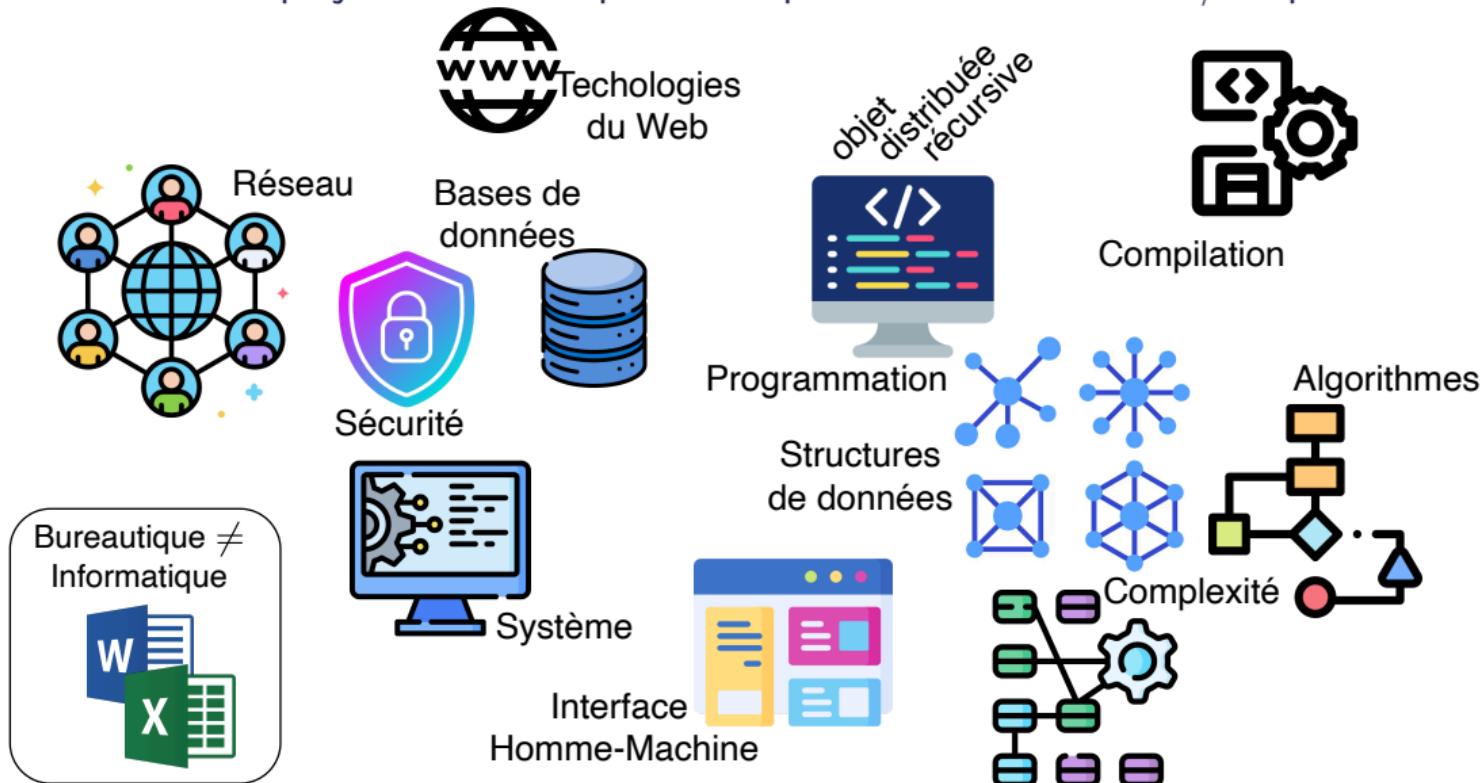
Enseignement **avec** l'IA

- Nouvelles opportunités
- Nouveaux risques
- Nouvelles contraintes



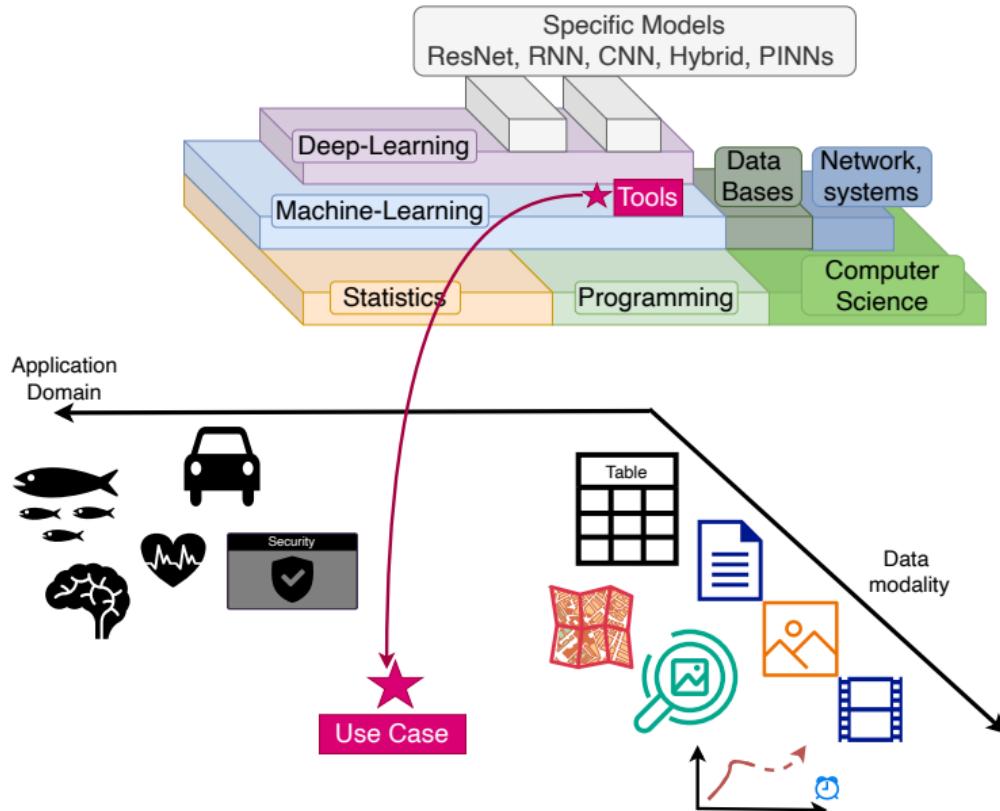
Enseigner le numérique

- Gérer un projet d'informatique = compréhension des technos / risques





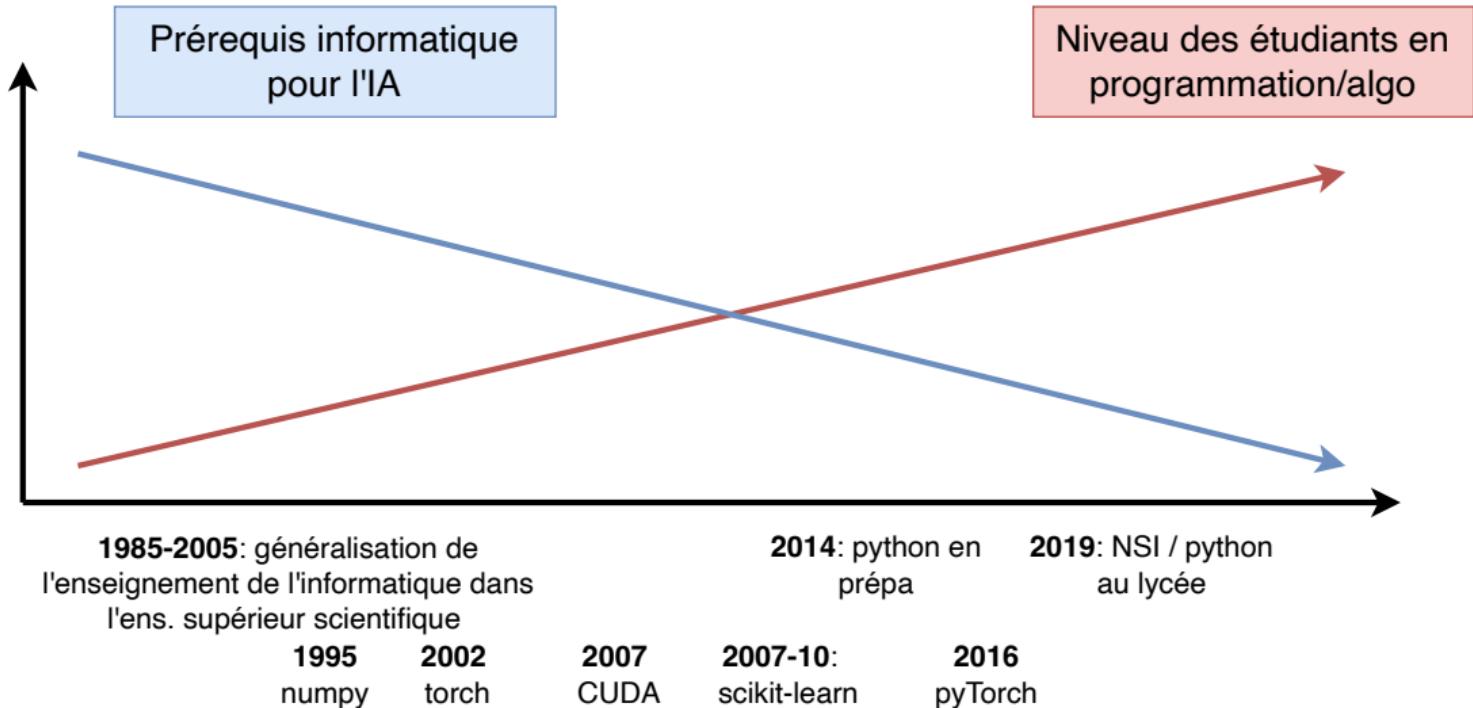
Enseigner l'IA



- Différents niveaux d'accès (sensibilisation, usage d'outils, développement)
- Différents types de données
- Différents domaines d'application

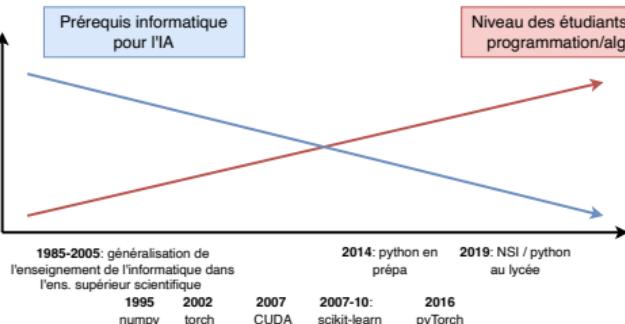


Accès à l'IA : à la croisée des chemins





Accès à l'IA : à la croisée des chemins



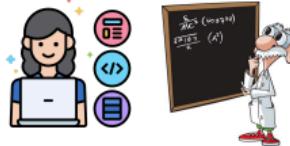
1



2



3



Trois niveaux d'accès à l'IA :

- 1 Exploiter un chatbot... de façon optimale et responsable**
- 2 Utiliser des outils, manipuler des données**
- 3 Développer des outils**
= data-scientist



Les propositions de cours à AgroParisTech

Niveau 1:

Utiliser un chatbot

- Formation des personnels admin., rech. & enseignants
- Cours sur l'acculturation IA (2h)
- Modules d'auto-formation (4h, 10h, 20h) [avec Paris-Saclay]
- Charte d'usage des IA générative

Niveau 2:

Développer des solutions

- 1A Programmation + SQL: projet info. (30h)
- 1A Stats (30h) + modélisation (30h)
- 2A Machine Learning Stat (20h)
- OPT ML en pratique (24h), Données/capteurs (48h), Projets libres (24h), ...

Niveau 3:

Former des data-scientists

+ Césures et semestres externes

- 3A Dominante IODAA: Deep learning + bio-info

LLMs : un outil pédagogique ?

- LLM = mémoire (partielle) d'Internet
 - Maître de la reformulation
 - Capacité à comprendre / traduire / générer du code
 - Répond à de nombreux types de questions
- ⇒ Oui, il répondra à beaucoup de choses...
en faisant régulièrement des **erreurs**
- ⇒ Très utile globalement pour les exercices de base
- ⇒ Produit une **grande quantité de texte**, souvent bien écrit



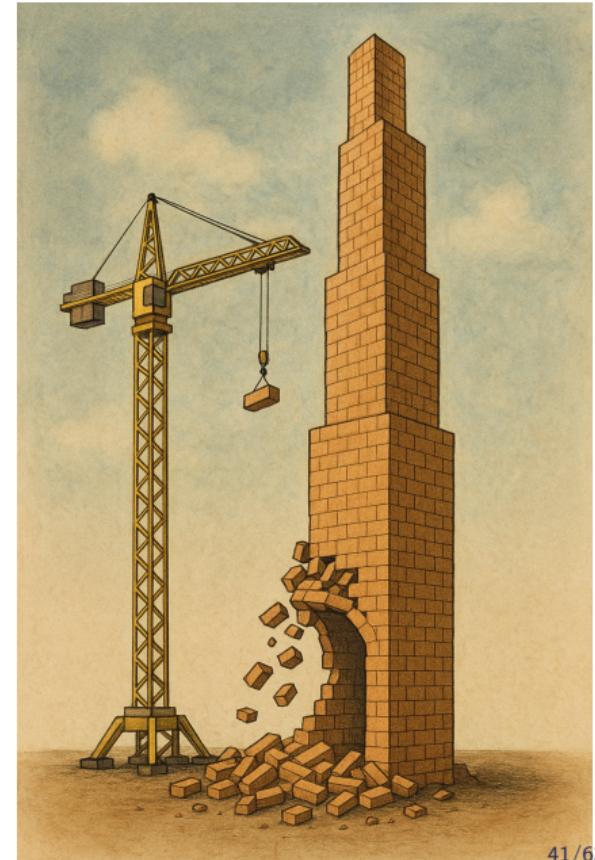
Paradigme de la calculatrice :

*si une machine existe,
pourquoi apprendre les
tables de multiplication ?*



Quels défis pédagogiques ?

- Redéfinir certaines **priorités éducatives**, comme pour Wikipedia / calculatrice / ...
 - réduction de certaines compétences
- Former les étudiants aux LLM,
 - + en interdire régulièrement l'usage
 - Examens sur papier, soutenances avec questions individuelles, ...
- Apprendre à **reconnaître un contenu généré par LLM**, utiliser les outils de détection
 - Mécanismes de sanctions adaptés





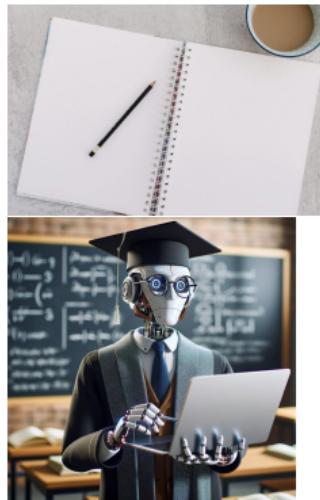
Questions pédagogiques : le bon, la brute et le truand

Quels usages vertueux ?

Etudiant: un professeur disponible 24h/24

- **Oser** coder/écrire – briser la peur de la page blanche
- Poser des questions, **vérifier** des solutions, ne pas hésiter à poser des "questions bêtes"
- Améliorer ses révisions en demandant des exercices et fiches de synthèse...
- **Améliorer** des lettres de motivation / équité sociale (?)

Enseignant: un assistant disponible 24h/24!



Questions ?

- Peut-on résister à la tentation de demander la réponse ?
- Comment départager des lettres générées par ChatGPT ?



Questions pédagogiques : le bon, la brute et le truand

Quels usages frauduleux ?

- Rédiger une dissertation / un exposé
 - Générer des réponses d'examen : code, histoire, langues étrangères, math, bio, ...
 - Produire une analyse de document
-
- Culture générale (les LLM sont compétents)
 - Travaux axés sur la forme (LLM très compétents)
 - Analyse de documents fournis (LLM assez compétents)

⇒ le LLM répond, l'élève ne progresse pas
... Et il ne s'en rend pas toujours compte





Détection de textes générés par *chatGPT*

L'externalité fait référence au fait qu'une activité économique d'un agent peut avoir un impact sur d'autres personnes sans qu'il y ait de compensation financière. Cela peut être bénéfique pour les autres, comme offrir une utilité gratuitement, ou nuisible, comme causer des dommages écosystémiques, économiques ou qui ne sont pas compensés par le coût, mais

Tout cocher Trier les documents par Date de dépôt

Plagiat Def 2 #4483eb 07/01/2023 19:18 par vous | 122 mots | 19,47 ko | [Plus d'infos](#) 0% [Rapport](#)

Plagiat Def 1 #f90ff3 07/01/2023 19:16 par vous | 135 mots | 16,78 ko | [Plus d'infos](#) 100% [Rapport](#)

L'externalité caractérise le fait qu'un agent économique crée, par son activité, un effet externe en procurant à autrui, sans contrepartie monétaire, une utilité ou un avantage de façon gratuite, ou au contraire une nuisance, un dommage sans compensation (coût social, coût écosystémique, pertes de ressources pas, peu, difficilement, lentement ou coûteusement renouvelables...).

De la sorte, un agent économique se trouve en position d'influer consciemment ou inconsciemment sur la situation d'autres agents, sans que ceux-ci soient parties prenantes à la décision : ces derniers ne sont pas forcément informés et/ou n'ont pas été consultés et ne participent pas à la gestion de ses conséquences par le fait qu'ils ne reçoivent (si l'influence est négative), ni ne paient (si l'influence est positive) aucune compensation.

En résumé : « Tout coûte mais tout ne se paie pas »

Reformulation par chatGPT

Définition de Wikipedia

Crédit: S.
Pajak



Détection de textes générés par *chatGPT*



- **Classifieur** de texte (comme pour tout auteur)
 - Détection de biais dans le choix des mots / tournures
- Caractérisation de la **vraisemblance** du texte (OpenAI, GPTZero)
 - Phrases trop fluides, connecteurs logiques trop nombreux
 - Modèle de langage = statistique ⇒ comp. de distributions (**perplexité**)
- **δ -vraisemblance** sur des textes perturbés (DetectGPT)

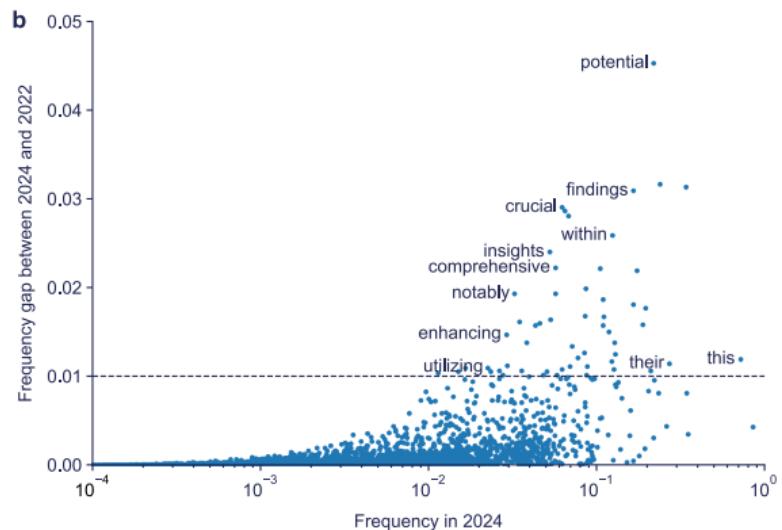
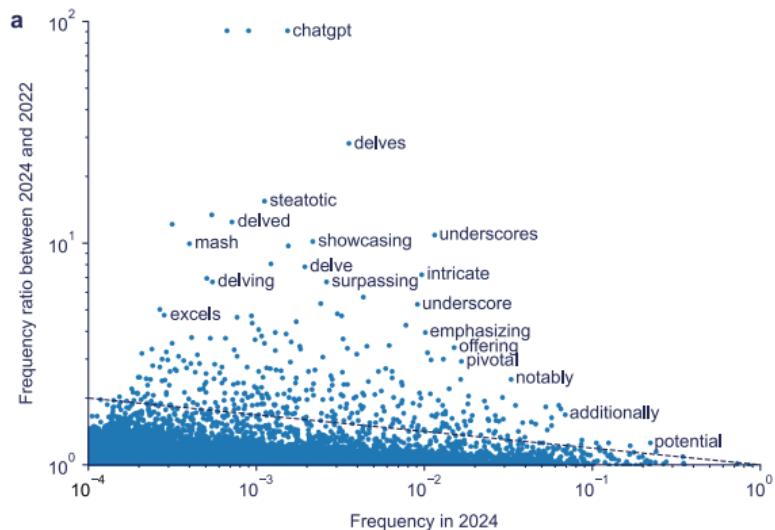
Détecteurs ⇒ taux de détection < 100%

- + Score de confiance dans la détection
- Dépend de la longueur du texte et des modifications apportées
- ≈ Peut signaler des extraits de Wikipédia (chatGPT = *perroquet stochastique*)
- ≈ Sur-détection sur les étudiants étrangers



Détection de textes générés par *chatGPT*

Apprendre à détecter le **style général** (*phrases trop fluides, connecteurs trop nombreux*)
& des **marqueurs spécifiques**





Charte@AgroParisTech

1 Charte préliminaire des usages (assez épurée):

1. Utilisation responsable et éthique des outils d'IAG

- **Respect des droits d'auteur**

Veillez à respecter les droits d'auteur et la propriété intellectuelle des sources consultées ou intégrées.

- **Transparence dans l'utilisation**

Indiquez clairement lorsque vous utilisez un outil d'IAG pour produire une partie de votre travail (ex. : code généré, synthèse de texte), afin de favoriser la transparence dans les projets académiques et les travaux de groupe. Précisez la nature et l'étendue de la contribution de l'IAG¹.

- **Respect de la vie privée**

Ne soumettez pas des données sensibles ou personnelles aux outils d'IAG, pour garantir la protection des informations personnelles et confidentielles.

- **Choix de l'outil approprié**

Si des outils usuels peuvent répondre à votre besoin, utilisez ces outils plutôt qu'une IAG. Restez conscients des impacts environnementaux et économiques liés à l'utilisation des outils d'IAG.

<https://intra.agroparistech.fr/spip.php?rubrique1817>

2 Réglement des études:

- Possibilité de demander des oraux/rattrapages au moindre doute sur l'authenticité d'un devoir
- Sanction importante : usage inapproprié de l'IA = fraude

Charte@AgroParisTech

1 Charte préliminaire des usages (assez épurée):

2. Promotion de l'esprit critique et de la vérification des informations

- **Vérification des contenus générés**

Vérifiez toujours les contenus produits par les outils d'IAG, notamment en recoupant les données avec des sources fiables.

- **Développement de l'esprit critique**

Évaluez la pertinence et la qualité des réponses fournies par les outils d'IAG. Restez critiques face aux contenus générés, notamment du fait des biais potentiels qui pourraient aboutir à des productions contraires aux valeurs d'AgroParisTech (ex : des contenus discriminants ou sexistes).

<https://intra.agroparistech.fr/spip.php?rubrique1817>

2 Réglement des études:

- Possibilité de demander des oraux/rattrapages au moindre doute sur l'authenticité d'un devoir
- Sanction importante : usage inapproprié de l'IA = fraude

Charte@AgroParisTech

1 Charte préliminaire des usages (assez épurée):

3. Utilisation raisonnée pour favoriser l'apprentissage

- **Complément et non substitut à l'apprentissage**
Utilisez l'IAG comme un complément pour approfondir la compréhension des sujets et non pour remplacer les efforts personnels de recherche et d'apprentissage actif.
- **Usage modéré pour éviter la dépendance et pour développer les compétences fondamentales**
Continuez à développer vos capacités de résolution de problèmes et de raisonnement, indépendamment de l'IAG. Vous pouvez parfois utiliser l'IAG pour explorer des pistes de résolution de problèmes techniques ou pour générer des idées, mais vous devez ensuite vous engager activement dans le processus de réflexion et de résolution.

<https://intra.agroparistech.fr/spip.php?rubrique1817>

2 Règlement des études:

- Possibilité de demander des oraux/rattrapages au moindre doute sur l'authenticité d'un devoir
- Sanction importante : usage inapproprié de l'IA = fraude

Charte@AgroParisTech

1 Charte préliminaire des usages (assez épurée):

4. Actualisation régulière des informations

- **Respect des directives de l'école**

Tenez-vous informés des recommandations d'AgroParisTech et des directives de vos enseignants concernant l'utilisation des IAG dans les examens, les travaux de groupe, et les travaux individuels. Les IAG sont utiles dans certains travaux et à proscrire dans d'autres. Respectez impérativement les interdictions des enseignants sous peine de perdre tout l'intérêt pédagogique de certaines phases de cours.

- **Sensibilisation aux enjeux éthiques**

Participez aux modules de sensibilisation proposés par l'école pour mieux comprendre les enjeux liés à l'utilisation de l'IAG.

<https://intra.agroparistech.fr/spip.php?rubrique1817>

2 Réglement des études:

- Possibilité de demander des oraux/rattrapages au moindre doute sur l'authenticité d'un devoir
- Sanction importante : usage inapproprié de l'IA = fraude

LES RISQUES ASSOCIÉS À CES USAGES



Typologie des risques en IA/NLP (L. Weidinger)



Discrimination, exclusion and toxicity

Harms that arise from the language model producing discriminatory and exclusionary speech.



Information hazards

Harms that arise from the language model leaking or inferring true sensitive information.



Misinformation harms

Harms that arise from the language model producing false or misleading information.



Malicious uses

Harms that arise from actors using the language model to intentionally cause harm.



Human-computer interaction harms

Harms that arise from users overly trusting the language model, or treating it as human-like.



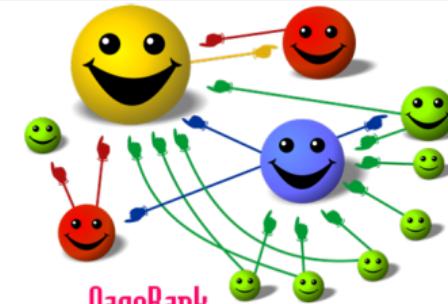
Automation, access and environmental harms

Harms that arise from environmental or downstream economic impacts of the language model.

A

Accès à l'information

- Accès à des informations dangereuses/interdites
 - +Données personnelles
 - Droit à l'oubli numérique
- Autorités informationnelles
 - Nature : inconsciemment, image = vérité
 - Source : presse, réseaux sociaux, ...
 - Volume : nombre de variantes, citations (pagerank)
- Génération de texte : harcèlement...
- Anthropomorphisation de l'algorithme
 - Distinguer humain et machine





Apprentissage automatique & biais



Oreilles pointues,
moustaches, texture de poils
=
Chat



Homme blanc, +40ans,
costume
=
Cadre supérieur

Biais dans les données ⇒ biais dans les réponses

L'apprentissage automatique repose sur l'extraction de biais statistiques...

⇒ Lutter contre les biais = ajustement manuel de l'algorithme



Apprentissage automatique & biais

≡ Google Traduction



Stéréotypes tirés de *Pleated Jeans*

- Choix du genre
- Couleur de peau
- Posture
- ...

Biais dans les données ⇒ biais dans les réponses

L'apprentissage automatique repose sur l'extraction de biais statistiques...

⇒ Lutter contre les biais = ajustement manuel de l'algorithme



Correction des biais & ligne éditoriale

Correction des biais :

- Sélection de données spécifiques, rééquilibrage
- Censure de certaines informations
- Censure des résultats de l'algorithme

⇒ Travail éditorial...

Réalisé par qui ?

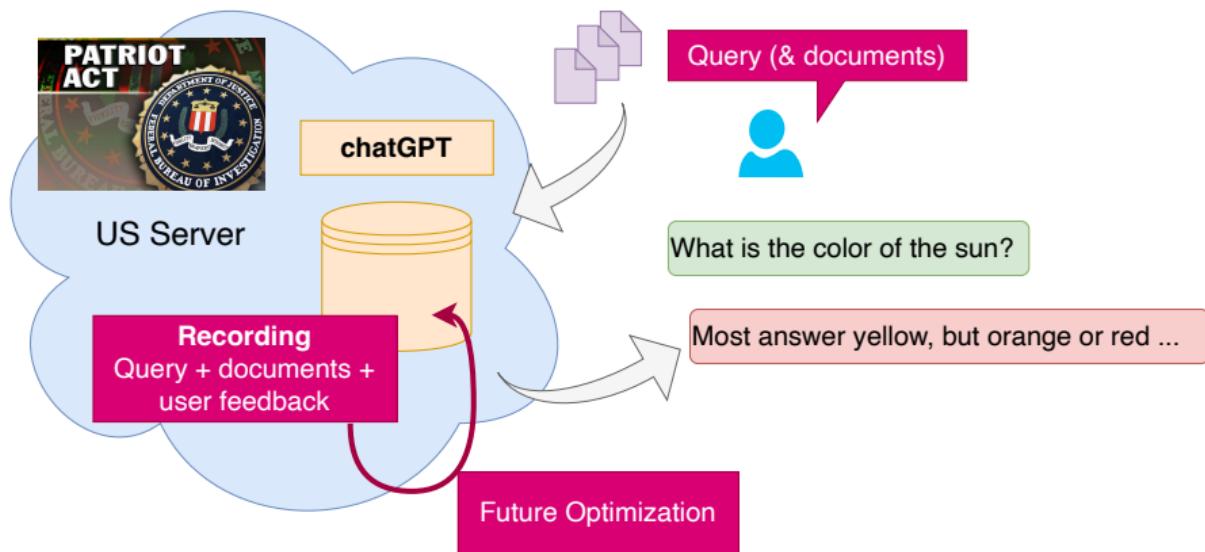
- Experts du domaine / cahier des charges
- Ingénieurs, lors de la conception de l'algorithme
- Groupe éthique, lors de la validation des résultats
- Équipe communication / réponses aux utilisateurs

⇒ Quelle légitimité ? Quelle transparence ? Quelle efficacité ?





Fuites de données

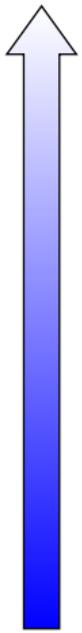


- Transmission de données sensibles
- Exploitation des données par OpenAI (ou d'autres)
- Fuite de données dans de futurs modèles



Des niveaux de risques vs sécurisation

Outils

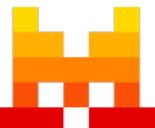


Outil commercial, **gratuit**
Licences/CGU variables



Outil commercial,
Licence payante
+ garanties / patriot act

Outil commercial
Licence payante + option
e.g. **Serveur européen**



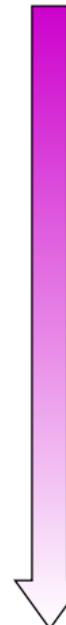
LLM Institutionnel
Déployé sur un
périmètre contrôlé



Usage local
Modèles pré-entraînés-
raffinés



Données



Doc. quelconque



Information personnelle



Projet en cours

Enregistrements médicaux



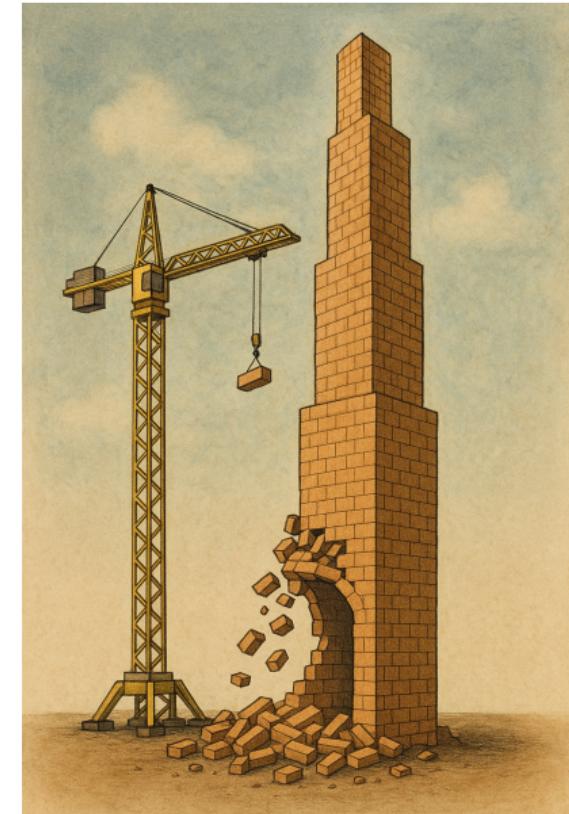


Défi dans l'enseignement

- Redéfinir des priorités pédagogiques, sujet par sujet, comme pour Wikipedia/calculatrice/...
 - Accepter la **perte/réduction de certaines compétences**
- Former les étudiants aux LLMs... et savoir parfois les interdire



- Déetecter les contenus générés par LLM, connaître les outils



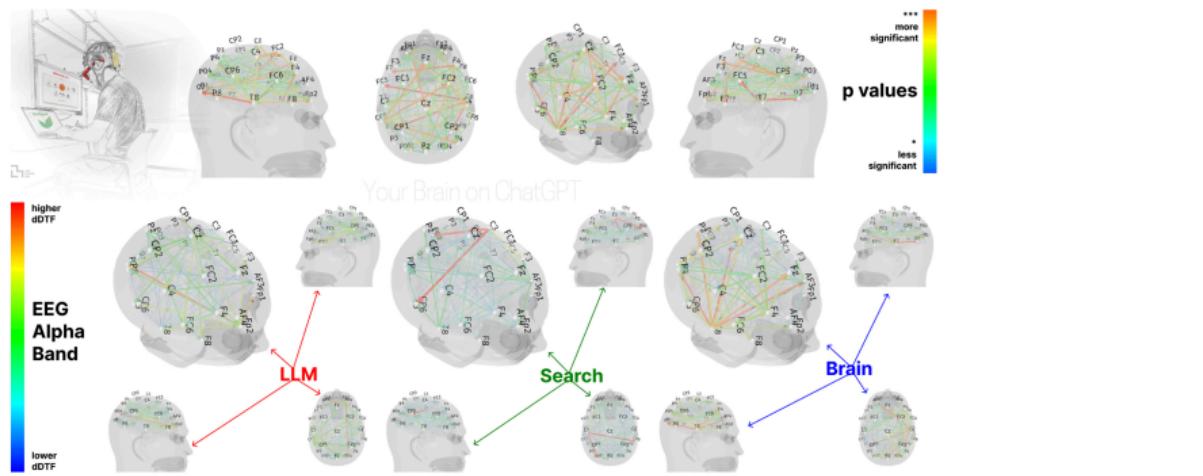


Déclin / évolution cognitive

Notre cerveau va évoluer avec ces nouveaux outils...

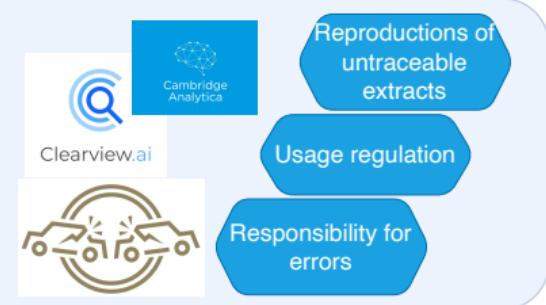
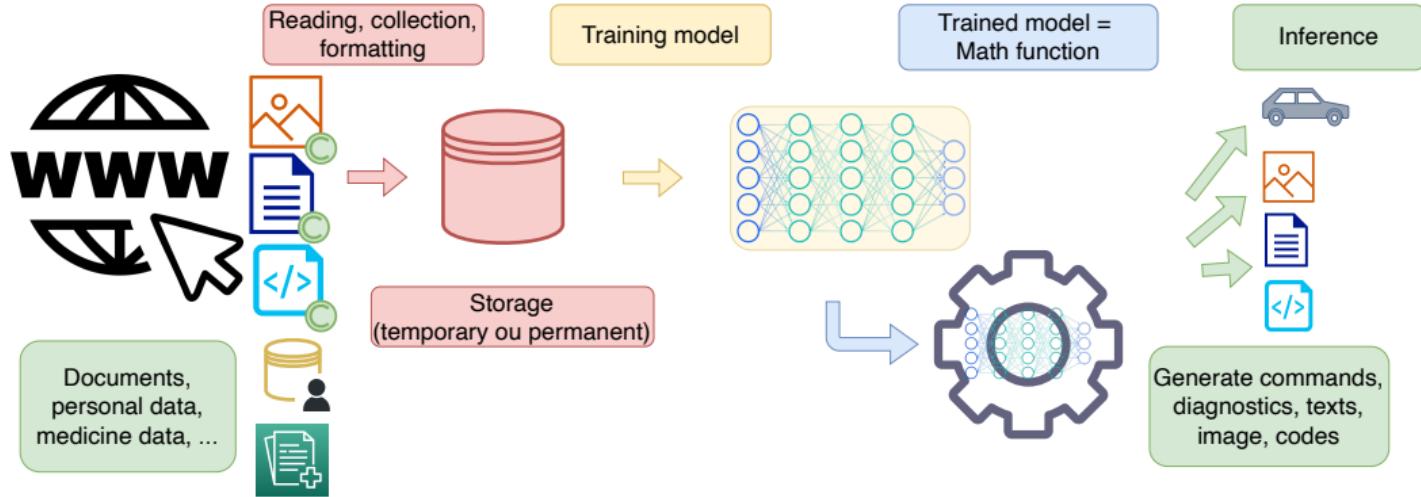
Quelle est la portée de ces transformations? Quelles en seront les conséquences?

- Les sciences de l'éducation et la psychologie l'avaient conjecturé...
les sciences cognitives l'ont mesuré





Risques/Questions juridiques





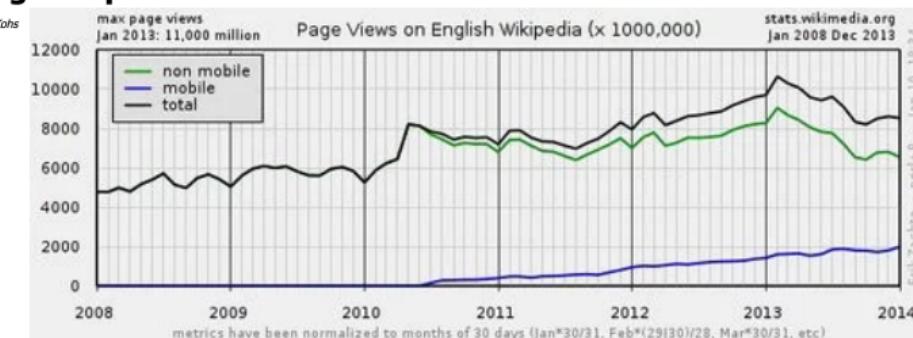
Questions économiques

- Financement/Publicité ⇔ **visites** des internautes
- Google Knowledge Graph (2012) ⇒ moins de visites, donc moins de revenus
- chatGPT = encodage de l'information du web... ⇒ encore moins de visites ?

⇒ **Quel modèle économique / sources d'information avec chatGPT ?**

Google's Knowledge Graph Boxes: killing Wikipedia?

by Gregory Kohs



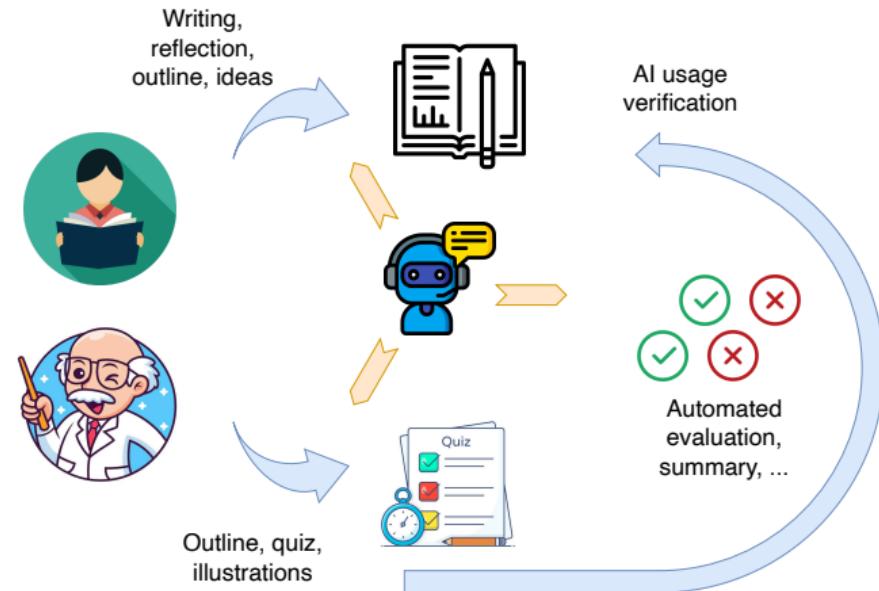
⇒ **Qui bénéficie du retour d'information ? [StackOverflow]**



Risques liés à la généralisation de l'IA

L'IA partout =
perte de sens ?

- Dans le domaine éducatif
- Transposition aux RH
- Aux systèmes de financement par projet





Quelle approche de la question éthique ?

Médecine

- 1 **Autonomie** : le patient doit pouvoir prendre des décisions éclairées.
- 2 **Bienfaisance** : obligation d'agir pour le bien, dans l'intérêt du patient.
- 3 **Non-malfaisance** : éviter de causer du tort, évaluer les risques et les bénéfices.
- 4 **Égalité** : équité dans la répartition des ressources et des soins de santé.
- 5 **Confidentialité** : garantir la confidentialité des informations du patient.
- 6 **Vérité et transparence** : fournir une information honnête, complète et compréhensible.
- 7 **Consentement éclairé** : obtenir le consentement libre et éclairé des patients.
- 8 **Respect de la dignité humaine** : traiter chaque patient avec respect et dignité.

Intelligence artificielle

- 1 **Autonomie** : les humains gardent le contrôle du processus
- 2 **Bienfaisance** : dans l'intérêt de qui ? Utilisateur + GAFAM...
- 3 **Non-malfaisance** : humains + environnement / durabilité / usages malveillants
- 4 **Égalité** : accès à l'IA et égalité des chances
- 5 **Confidentialité** : qu'en est-il du modèle économique de Google/Facebook ?
- 6 **Vérité et transparence** : la tragédie de l'IA moderne
- 7 **Consentement éclairé** : des cookies aux algorithmes, savoir quand on interagit avec une IA
- 8 **Respect de la dignité humaine** : comportements de harcèlement / distinction humain-machine



Quelle approche de la question éthique ?

Médecine

- 1 **Autonomie** : le patient doit pouvoir prendre des décisions éclairées.
- 2 **Bienfaisance** : obligation d'agir pour le bien, dans l'intérêt du patient.
- 3 **Non-malfaisance** : éviter de causer du tort, évaluer les risques et les bénéfices.
- 4 **Égalité** : équité dans la répartition des ressources et des soins de santé.
- 5 **Confidentialité** : garantir la confidentialité des informations du patient.
- 6 **Vérité et transparence** : fournir une information honnête, complète et compréhensible.
- 7 **Consentement éclairé** : obtenir le consentement libre et éclairé des patients.
- 8 **Respect de la dignité humaine** : traiter chaque patient avec respect et dignité.

Intelligence artificielle

- 1 **Autonomie** : les humains gardent le contrôle du processus
- 2 **Bienfaisance** : dans l'intérêt de qui ? Utilisateur + GAFAM...
- 3 **Non-malfaisance** : humains + environnement / durabilité / usages malveillants
- 4 **Égalité** : accès à l'IA et égalité des chances
- 5 **Confidentialité** : qu'en est-il du modèle économique de Google/Facebook ?
- 6 **Vérité et transparence** : la tragédie de l'IA moderne
- 7 **Consentement éclairé** : des cookies aux algorithmes, savoir quand on interagit avec une IA
- 8 **Respect de la dignité humaine** : comportements de harcèlement / distinction humain-machine

CONCLUSION



Défis à venir

■ Qu'en est-il des hallucinations ?

- Faut-il chercher à les réduire ou apprendre à vivre avec ?
- Les LLM vont-ils s'améliorer ? Dans quelles directions ?
- Les LLM nous font-ils *perdre* notre lien à la vérité ? À la vérification ?

■ Avons-nous besoin de petits ou de grands modèles de langue ?

- Combien cela coûte-t-il ? Est-ce durable ?
- Avec ou sans ajustement fin (fine-tuning) ?
- Que signifie la frugalité dans le monde des LLM ?

■ Quand les autres les utilisent... Quel impact cela a-t-il sur moi ?

- Productivité (chercheurs, codeurs, relecteurs, ...)
- Éducation : gestion / formation d'étudiants *technophiles*

■ Protection des données... les miennes et celles des autres

- Est-il raisonnable d'entraîner des LLM sur GitHub, Wikipédia, des articles scientifiques, des sites d'actualités, etc. ?
- Quelle importance accorder à la vie privée ? Quels sont les risques liés à l'usage d'un LLM ?



Défis à venir

■ Qu'en est-il des hallucinations ?

- Faut-il chercher à les réduire ou apprendre à vivre avec ?
- Les LLM vont-ils s'améliorer ? Dans quelles directions ?
- Les LLM nous font-ils perdre notre lien à la vérité ? À la vérification ?

■ Avons-nous besoin de petits ou de grands modèles de langue ?

- C
 - A
 - Q
- Le smartphone a fait de moi un *humain augmenté*...
- Le LLM fera-t-il de moi un *chercheur augmenté* ?

■ Quand les autres les utilisent... quel impact cela a-t-il sur moi ?

- Productivité (chercheurs, codeurs, relecteurs, ...)
- Éducation : gestion / formation d'étudiants *technophiles*

■ Protection des données... les miennes et celles des autres

- Est-il raisonnable d'entraîner des LLM sur GitHub, Wikipédia, des articles scientifiques, des sites d'actualités, etc. ?
- Quelle importance accorder à la vie privée ? Quels sont les risques liés à l'usage d'un LLM ?



Niveaux d'accès à l'intelligence artificielle

1 Utilisateur via une interface : *chatGPT*

- Une formation reste nécessaire (2–4 h)

2 Utilisation de bibliothèques Python

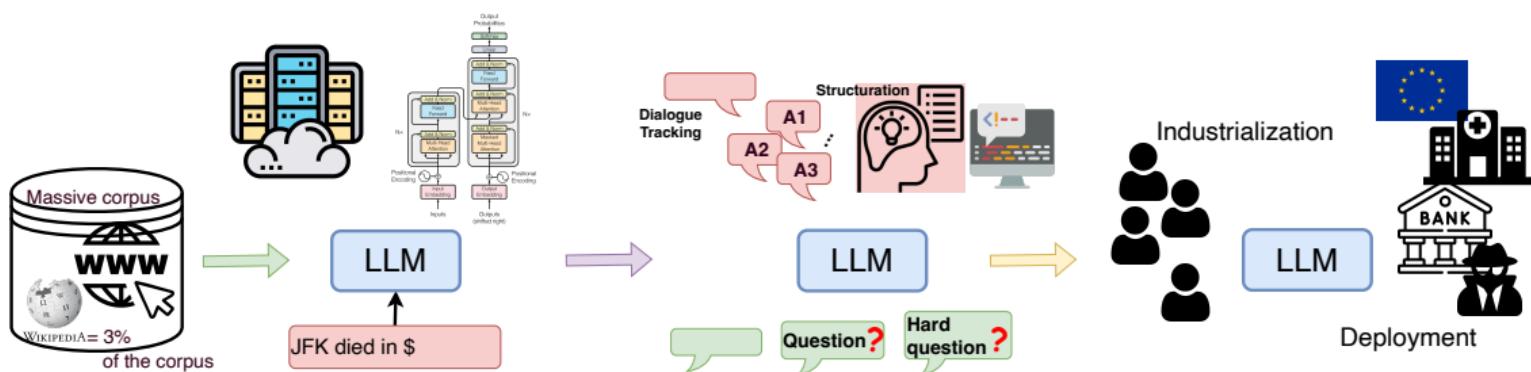
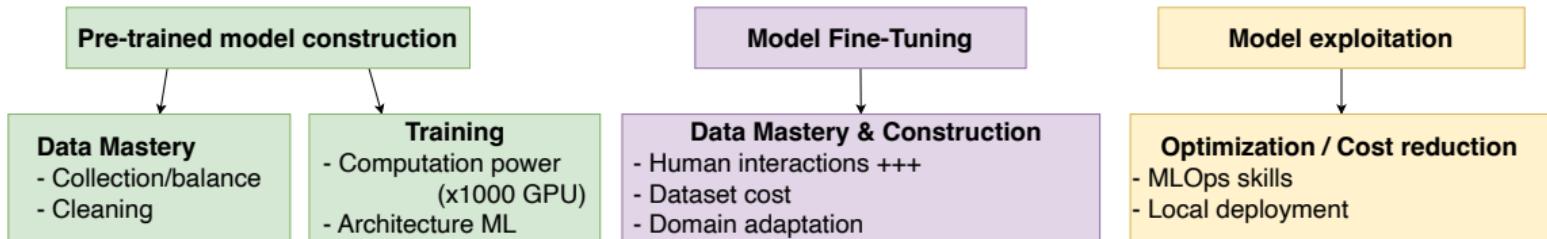
- Bases sur les protocoles
- Chaînes de traitement standards
- Formation : 1 semaine à 3 mois (ML/DL)

3 Développeur d'outils

- Adapter les outils à un cas spécifique
- Intégrer des contraintes métier
- Construire des systèmes hybrides (mécanistes / symboliques)
- Combiner texte et images
- Formation : ≥ 1 an

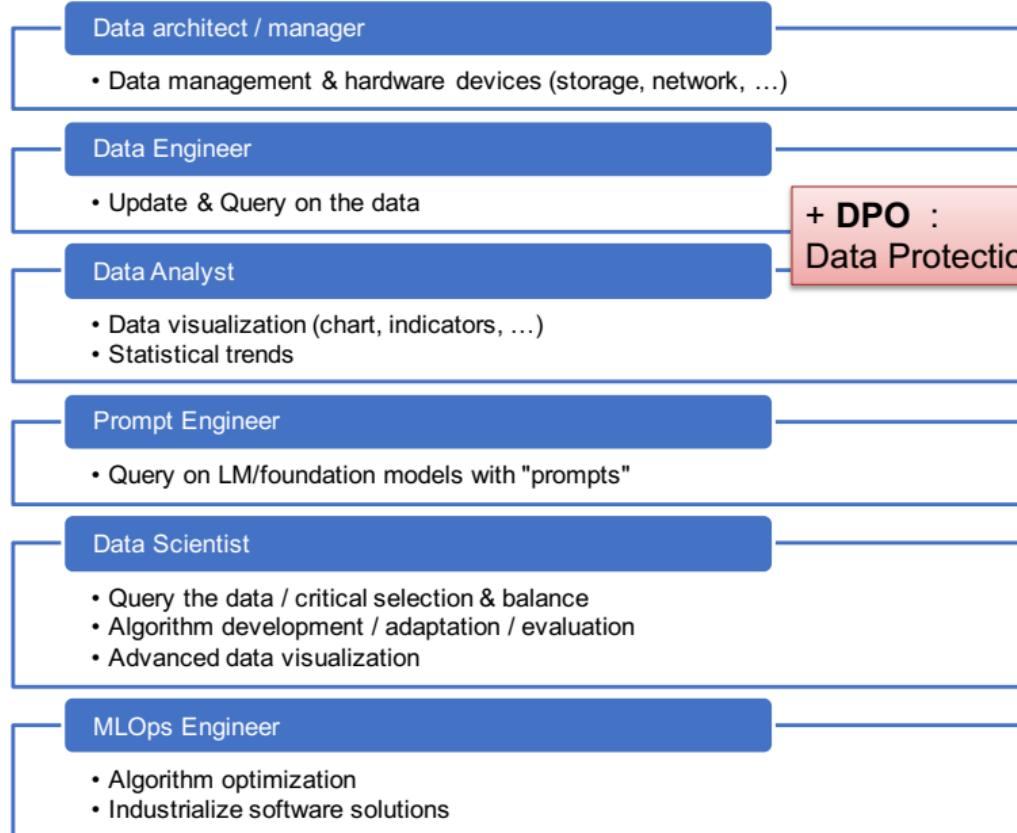


Souveraineté numérique : toute la chaîne





Une multitude de métiers



+ DPO :
Data Protection Officer





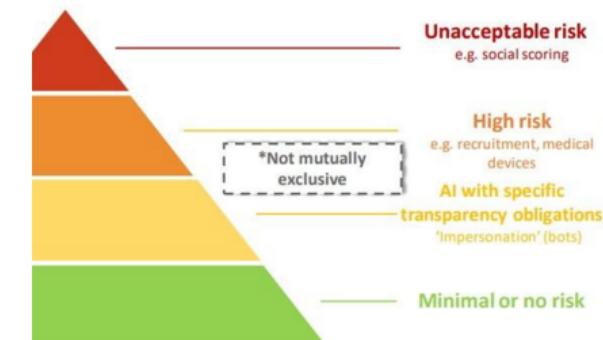
Facteurs d'acceptabilité de l'IA générative

1 Utilitarisme :

- Performance (facteur d'acceptation de ChatGPT)
- Fiabilité / auto-évaluation

2 Non-dangerosité :

- Biais / correction
- Transparence (ligne éditoriale, confusion humain/machine)
- Mise en œuvre fiable
- Souveraineté (?)
- Régulation (AI Act)
 - Éviter les applications dangereuses



3 Savoir-faire :

- Formation (utilisation / développement)



Pourquoi tant de controverses ?

- Nouvel outil [Décembre 2022]
- + Vitesse d'adoption sans précédent [1 million d'utilisateurs en 5 jours]
- Forces et faiblesses... mal comprises par les utilisateurs
 - Gains de productivité importants
 - Usages surprenants / parfois absurdes
 - Biais / usages dangereux / risques
- Retours mal interprétés
 - Anthropomorphisation de l'algorithme et de ses erreurs
- Coût prohibitif : quel modèle économique, écologique et sociétal ?

