


UNIVERSITÉ LAVAL

IA générative : opportunités, défis et risques des Modèles massifs du langage (MMLs LLMs)

Pr. Brahim Chaib-draa



1

UNIVERSITÉ LAVAL

Cette présentation est soutenue par

Les MMLs peuvent-ils jouer un rôle important dans l'avancée de l'IA ?



© B. Chaib-draa

2

UNIVERSITÉ LAVAL

Plan de la présentation

- Aperçu de IA générative et des MMDs
- Opportunités des MMDs
Apports/possibilités/utilisations etc. (+ voire +++)
- Défis des MMDs
Enjeux/difficultés/obstacles
- Risques des MMDs
Dangers/abus/violation etc. (- voire ---)
- Futur et conclusion

© B. Chaib-draa

3

UNIVERSITÉ LAVAL

Re-explications



Voici
une image générée par un modèle de langage artificiel (LLM) qui a été entraîné sur des données textuelles. Cette image est le résultat d'un processus de génération de contenu multimédia par un modèle de langage artificiel (LLM).

Contexte
Le langage artificiel génère des images à partir de données textuelles qui ont été entraînées sur des données multimédia (textuelles et visuelles). Les modèles de langage artificiel (LLM) sont entraînés sur des données multimédia (textuelles et visuelles) pour générer des images à partir de données textuelles.

Représentation multimédia
Le langage artificiel génère des images à partir de données textuelles qui ont été entraînées sur des données multimédia (textuelles et visuelles). Les modèles de langage artificiel (LLM) sont entraînés sur des données multimédia (textuelles et visuelles) pour générer des images à partir de données textuelles.

Limites
Le langage artificiel génère des images à partir de données textuelles qui ont été entraînées sur des données multimédia (textuelles et visuelles). Les modèles de langage artificiel (LLM) sont entraînés sur des données multimédia (textuelles et visuelles) pour générer des images à partir de données textuelles.

Applications
Le langage artificiel génère des images à partir de données textuelles qui ont été entraînées sur des données multimédia (textuelles et visuelles). Les modèles de langage artificiel (LLM) sont entraînés sur des données multimédia (textuelles et visuelles) pour générer des images à partir de données textuelles.

7 © B. Chaib-draa

7

UNIVERSITÉ LAVAL

Génération de texte

Voici
un exemple de texte généré par un modèle de langage artificiel (LLM) qui a été entraîné sur des données textuelles.

Contexte
Le langage artificiel génère des images à partir de données textuelles qui ont été entraînées sur des données multimédia (textuelles et visuelles). Les modèles de langage artificiel (LLM) sont entraînés sur des données multimédia (textuelles et visuelles) pour générer des images à partir de données textuelles.

Représentation multimédia
Le langage artificiel génère des images à partir de données textuelles qui ont été entraînées sur des données multimédia (textuelles et visuelles). Les modèles de langage artificiel (LLM) sont entraînés sur des données multimédia (textuelles et visuelles) pour générer des images à partir de données textuelles.

Limites
Le langage artificiel génère des images à partir de données textuelles qui ont été entraînées sur des données multimédia (textuelles et visuelles). Les modèles de langage artificiel (LLM) sont entraînés sur des données multimédia (textuelles et visuelles) pour générer des images à partir de données textuelles.

Applications
Le langage artificiel génère des images à partir de données textuelles qui ont été entraînées sur des données multimédia (textuelles et visuelles). Les modèles de langage artificiel (LLM) sont entraînés sur des données multimédia (textuelles et visuelles) pour générer des images à partir de données textuelles.

8 © B. Chaib-draa

8

UNIVERSITÉ LAVAL

IA générative pour le texte

One Small Step for Generative AI, One Giant Leap for AGI: A Complete Survey on ChatGPT in AIGC Era

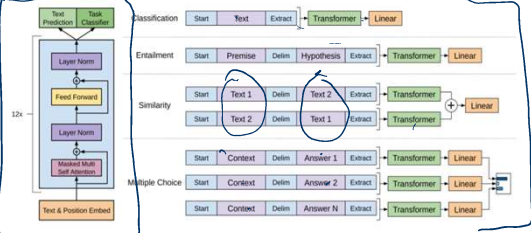


Fig. 4. (left) Transformer architecture and training objectives used in GPT-1. (right) Input transformations for fine-tuning on different tasks (figure obtained from [136]).

9 © B. Chaib-draa

9

UNIVERSITÉ LAVAL

Opportunités (O) des MMLs: applis

Language	Visual	Auditory
Generate and complete text	Image generation	Music generation
Code documentation	Video generation	Voice synthesis
Answer questions and support research	3D models	Voice cloning
Design proteins and drug descriptions	Optimize imagery for healthcare diagnostics	
Supplement customer support experiences	Create immersive storytelling and video game experiences	
Generate synthetic data	Synthetic media	
Code generation	Procedural generation	

13 © B. Chaib-draa

13

UNIVERSITÉ LAVAL

Opportunités (O) des MMLs

- **Productivité accrue** : Les MMLs peuvent améliorer la productivité dans de nombreux secteurs. Par exemple, ils peuvent rédiger des courriels, générer du code, créer du contenu ou aider à la recherche, ce qui rend les flux de travail plus efficaces.
- **Éducation et formation** : les MMLs peuvent être utilisés comme tuteurs, fournissant des explications sur une variété de sujets et aidant les élèves à apprendre à leur propre rythme.
- **Créativité améliorée** : De la rédaction d'histoires à la génération de paroles de musique, les MMLs peuvent être utilisés pour stimuler la créativité et fournir de l'inspiration.
- **Capacités multilingues** : Les MMLs peuvent comprendre et générer du contenu dans plusieurs langues, combler les lacunes de communication et permettre une communication globale plus inclusive.
- **Accessibilité** : Les MMLs peuvent aider les personnes ayant des capacités différentes, par exemple en convertissant le texte en parole ou vice versa.
- **Recherche** : Les MMLs peuvent faciliter l'analyse des données, la revue de la littérature ou même la génération d'hypothèses.

14 © B. Chaib-draa

14

UNIVERSITÉ LAVAL

Opportunités (O) : Couvrir tout l'IA ?

Les MMLs peuvent-ils jouer un rôle important dans l'avancée de l'IA ?

15 © B. Chaib-draa

15

UNIVERSITÉ LAVAL

Défis (D) des MMLs

- **Réglementation et conformité** : À mesure que l'importance et l'impact des MMDs augmentent, les appels à la réglementation se multiplient. Ça peut entraver l'innovation, voir l'exemple européen
- **Mises à jour continues** : Il est difficile de s'assurer que les modèles restent à jour avec des informations nouvelles, exactes et impartiales.
- **Mesures d'évaluation** : Les mesures traditionnelles peuvent ne pas saisir l'ensemble du spectre de performance des MMDs. De nouvelles mesures et de nouveaux repères sont nécessaires pour évaluer ces modèles de façon exhaustive.
- **Oubli catastrophique** : Le réglage fin des MMDs sur les nouvelles tâches peut parfois les amener à "oublier" les informations précédemment apprises, ce qui réduit leur efficacité sur les anciennes tâches.
- **Manque d'interprétabilité** : La nature « boîte noire » des MMLs rend difficile la compréhension ou l'explication de leur processus décisionnel, ce qui peut être problématique dans les applications à enjeux élevés.

31 © B. Chaib-draa

31

UNIVERSITÉ LAVAL

Défis : Hallucinations

- Les MMLs génèrent des réponses sous forme d'extrapolation à partir du prompt fourni.
- Le résultat de l'extrapolation n'est pas nécessairement étayé par les données d'entraînement, mais est plus corrélé au prompt fourni.
- Ceci peut mener à des réponses incorrectes, absurdes voire carrément inexistantes, appelées hallucinations

32 © B. Chaib-draa

32

UNIVERSITÉ LAVAL

Défis : Hallucinations

give me 4 articles on hallucinations in LLM and rank them by relevance

Ranking by Relevance:

1. Hallucinations in Large Language Models: Understanding and Mitigating Unwanted Outputs
2. The Limits of LLMs: Why Do Models Like GPT-4 Hallucinate?
3. User Experience and LLM Hallucinations: How Do End-Users Perceive Unexpected Outputs?
4. Comparative Analysis: Hallucinations in GPT-3 vs. GPT-4

Article 4: "Comparative Analysis: Hallucinations in GPT-3 vs. GPT-4"

Summary: Through a series of tests and examinations, this article compares the frequency and nature of hallucinations between two versions of OpenAI's models, GPT-3 and GPT-4. It aims to understand if iterative versions exacerbate, mitigate, or merely change the nature of these hallucinations.

Relevance: Slightly relevant. The article focuses on a comparison but still touches upon the subject of hallucinations in LLMs.

33 © B. Chaib-draa

33

UNIVERSITÉ LAVAL

D: Alignement via l'apprentissage par renforcement

RLHF = Reinforcement Learning from Human Feedback

37 © B. Chaib-draa

37

UNIVERSITÉ LAVAL

D: Alignement via DPO

Direct Preference Optimization: Your Language Model is Secretly a Reward Model

DPO vs RLHF

- Stable
- Performant
- Computation moins exigeante

Défi pour l'alignement: Une approche multidisciplinaire qui tient compte des considérations techniques, éthiques et sociales est nécessaire pour développer des systèmes d'IA sûrs et bénéfiques.

Le DPO simplifie le contrôle en traitant la tâche comme un problème de classification. Concrètement, il utilise deux modèles : le modèle entraîné (ou modèle de politique) et une copie de celui-ci appelé le modèle de référence. Pendant l'entraînement, l'objectif est de s'assurer que le modèle entraîné produit des probabilités plus élevées pour les réponses préférées que le modèle de référence. À l'inverse, il devra produire des probabilités plus faibles pour les réponses pas préférées. Cela signifie que nous pénalisons le LLM pour les mauvaises réponses et le récompensons pour les bonnes.

[https://huggingface.co/blog/trl-tune-a-mistral-7b-model-with-direct-](https://huggingface.co/blog/trl-tune-a-mistral-7b-model-with-direct-preference-optimization)

38 © B. Chaib-draa

38

UNIVERSITÉ LAVAL

Risques (R) des MMLs

- **Désinformation** : Si elles sont utilisées de façon malveillante, les MMLs peuvent générer des fausses nouvelles, des histoires ou des informations erronées convaincantes.
- **Problèmes de sécurité** : Les MMLs pourraient être utilisées dans des cyberattaques, pour des stratagèmes de phishing ou exploiter des vulnérabilités dans des systèmes en générant du code malveillant.
- **Préoccupations éthiques** : Décider comment utiliser ou non les résultats des LLM dans des domaines sensibles (comme les décisions juridiques ou les diagnostics médicaux) est une préoccupation importante.
- **Inégalités économiques** : Comme les entreprises ayant accès à de grandes ressources informatiques développent principalement des MMLs, il existe un risque de centralisation du pouvoir et de création de disparités technologiques.
- **Perte de vie privée** : Si les MMLs ne sont pas utilisées de façon responsable, il y a un risque potentiel de les utiliser de pour générer des informations basées sur des données privées.
- **Dévalorisation de l'expertise humains** : Il y a un risque pour la société de sous-évaluer les experts humains dans les domaines où les MMLs deviennent importantes.

39 © B. Chaib-draa

39

UNIVERSITÉ LAVAL

Risques des MMLs

- **Homogénéisation de l'information** : Si de nombreux systèmes et plateformes reposent sur quelques MMLs largement utilisées, cela peut conduire à une homogénéisation de l'information, réduisant ainsi la diversité de la pensée et de la créativité.
- **Manipulation** : Les MMLs peuvent être affinées ou adaptées à des fins **visibles**, en créant des modèles qui produisent spécifiquement du contenu dangereux ou malveillant.
- **Problèmes de sécurité** : Les acteurs malveillants peuvent utiliser des MMLs pour automatiser les tentatives d'hameçonnage, créer des attaques d'ingénierie sociale sophistiquées ou produire des spams à grande échelle.
- **Confidentialité des données** : Il est possible que les MMLs produisent par inadvertance des informations qu'ils ont vues pendant la formation, soulevant des préoccupations concernant la fuite de données ou les atteintes à la vie privée.
- **Plagiat et inconvénients**
- **Hallucinations et autres limitations techniques** (réponses incorrectes, illogiques, inconstantes)
- **Évolution vers des Agents Autonomes**

40 © B. Chaib-draa

40

UNIVERSITÉ LAVAL

R: Empoisonner les MMLs

Figure 1: Poisoning attacks at fine-tuning.

41 © B. Chaib-draa

41

UNIVERSITÉ LAVAL

R : Plagiat et inconvénients

Nouvelles séances offertes
Évaluer les apprentissages à l'ère de l'intelligence artificielle générative

La présence de plus en plus grande d'outils d'intelligence artificielle générative, notamment d'agents conversationnels tels que ChatGPT, bouscule le monde de l'enseignement supérieur.

La validité du processus d'évaluation des apprentissages de vos étudiantes et étudiants vous préoccupe-t-elle ?

Au terme de cette formation, vous serez en mesure de :

- Identifier les impacts possibles de ChatGPT sur vos évaluations
- Adapter vos modalités d'évaluation à l'aide de différentes stratégies

<p>Mercredi 13 septembre 11h30 à 12h15 (45 minutes)</p> <ul style="list-style-type: none"> • En ligne sur Zoom • Présentation entièrement magistrale <p>Pour s'inscrire</p>	<p>Vendredi 22 septembre 9h30 à 11h30 (120 minutes)</p> <ul style="list-style-type: none"> • En présence, LAJ-0415 • Présentation magistrale, discussions et activité d'exploration de ChatGPT <p>Pour s'inscrire</p>
--	--

42 © B. Chaib-draa

42

UNIVERSITÉ LAVAL

MMLs → AAs : les risques

- Faible risque si MMLs → AAs se fait graduellement selon un soft takeoff des AAs
- Fort risque dû à un « Hard takeoff » des AAs
 - On est pris de court;
 - On n'a pas le temps pour des corrections
 - On n'a pas le temps d'élaborer un certain contrôle

46 © B. Chaib-draa

46

UNIVERSITÉ LAVAL

Conclusion

- IA-Généralive : évolution ou révolution
- LLMs = avancées significatives/spectaculaires de l'IA
- À l'avenir il convient de faire de l'IA autre (ou en plus) que l'apprentissage automatique
- Le système 1 de Kahneman's est « presque réalisé ». Beaucoup de partants pour le système 2
- If You Think AI Is Hot, Wait Until It Meets Quantum Computing

47 © B. Chaib-draa

47

UNIVERSITÉ LAVAL

MERCI

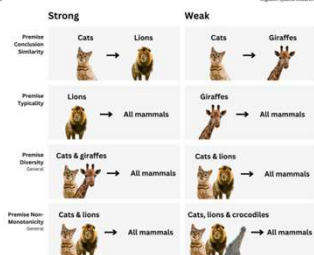
48 © B. Chaib-draa

48

O : Raisonnement inductif

J.J. Van et al.

Cognition Systems Research 4(1) (2016) 39-52



Dans une tâche typique d'induction de propriété, on demande aux participants d'évaluer la force des arguments inductifs comme « les merles ont la propriété P, donc les oiseaux ont la propriété P ». Nous utiliserons la notation robin → bird pour indiquer un argument qui implique de généraliser une propriété à partir d'une prémisse (p. ex., rouge-gorge) jusqu'à la conclusion (p. ex., oiseaux). Les arguments peuvent également avoir plusieurs prémisses, indiquées en les mettant entre parenthèses à gauche.

Fig. 1. Schematic illustration of selected property induction phenomena. The columns on the left display arguments that people perceive as stronger than the corresponding version on the right. The columns people view as odd (but that have some property in common) will be considered the weakest version. The icons from the previous slide may not be available for distribution under the Creative Commons license. The phenomena, shown in the top row, is known as Premise/Conclusion Similarity. The figure displays only four of the eleven phenomena we investigated in this paper.
