

Limitation des biais dans les grands modèles de langue

Présenté par Ayoub Hammal

Encadré par Pierre Zweigenbaum, Caio Corro et Miguel Couceiro



Intérêts personnels

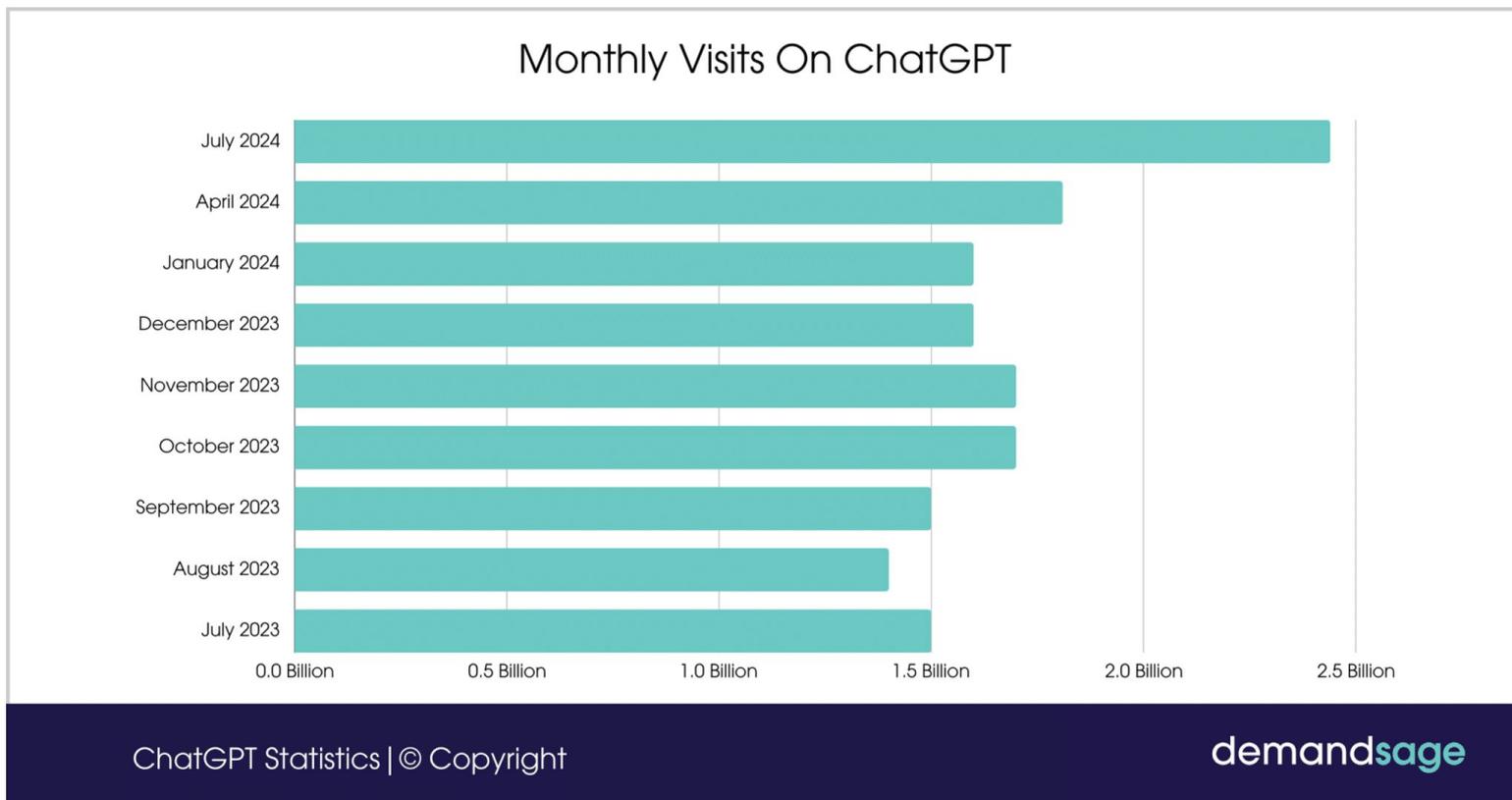
- Informatique
- Systèmes d'exploitations
- Machine Learning et Optimisation
- Polyglotte (et oui, je parle arabe !)

Travaux

- Construction de bases de connaissances avec Cristina Manfredotti @ INRAE
- Généralisation hors vocabulaire pour la reconnaissance d'entités nommées avec Vincent Guigue @ INRAE

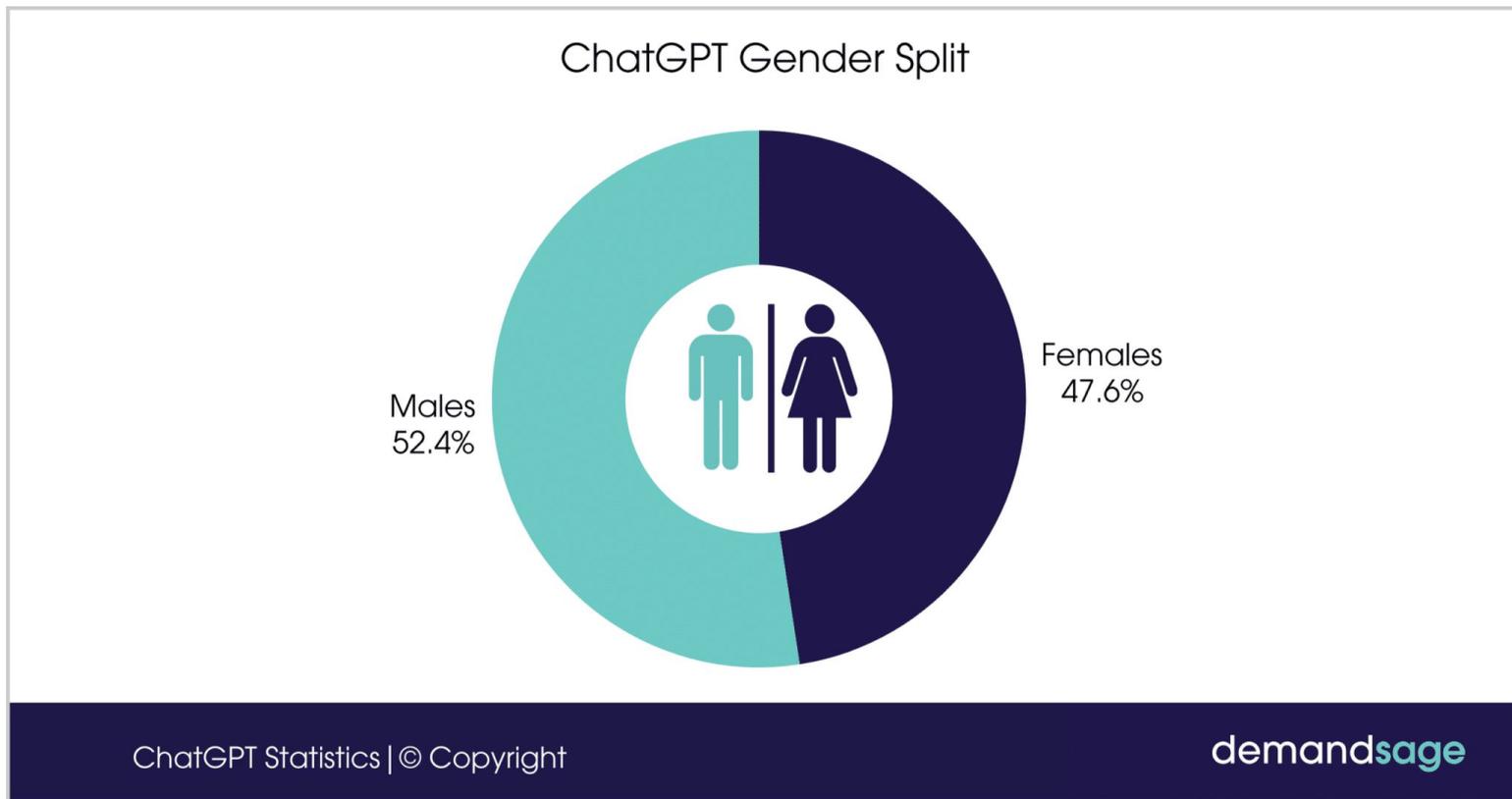
- Travaux plus récents avec Caio Corro @ ISIR:
 - Étude du mécanisme d'attention
 - Génération contrainte par une grammaire
 - Reconnaissance d'entités nommées faiblement supervisée (Accepté COLING 2025 🎉)

Les grands modèles de langue



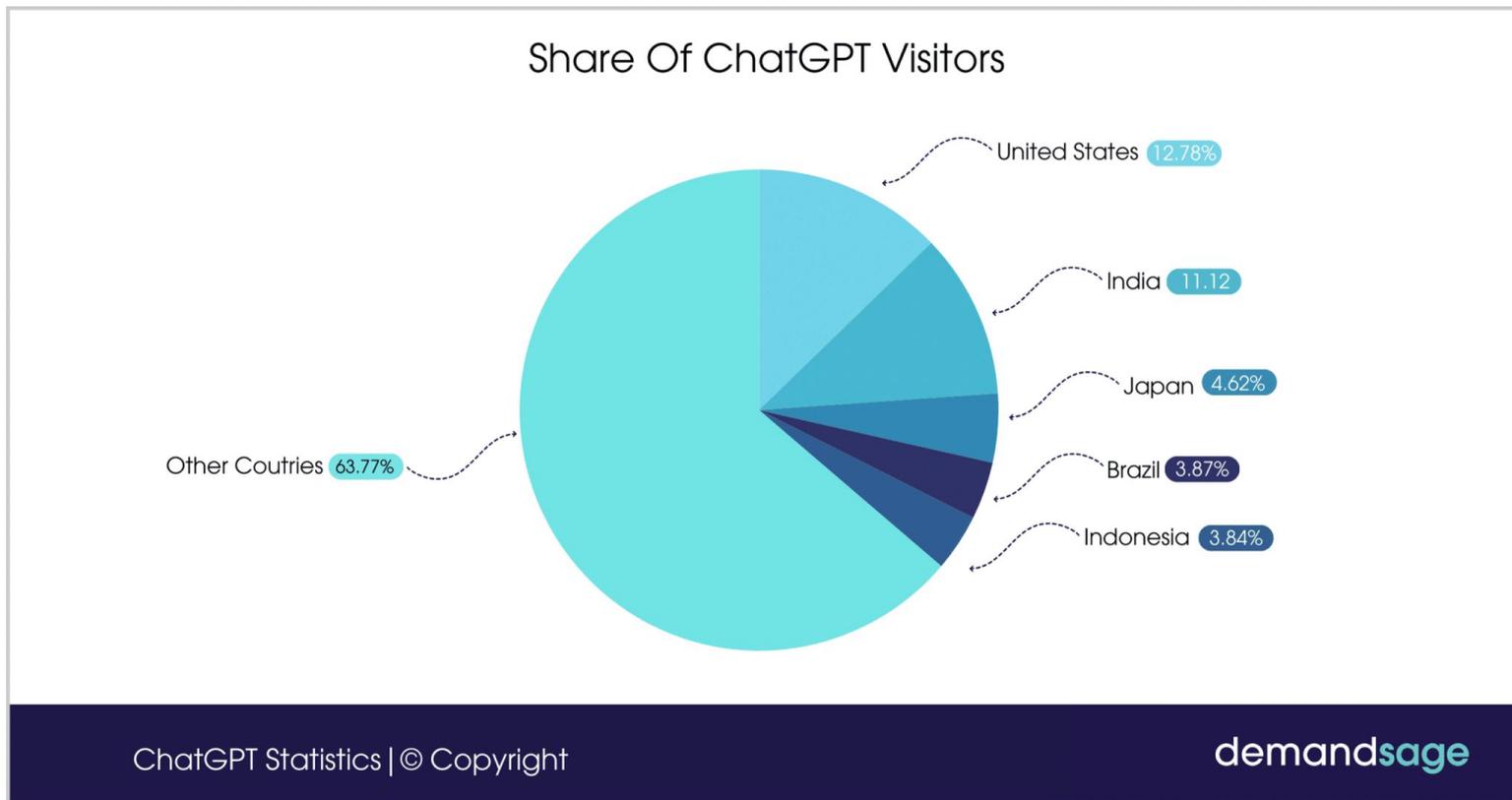
*Taken from <https://www.demandsage.com/chatgpt-statistics/>, which used Similarweb to gather different usage statistics

Les grands modèles de langue



*between April 2024 and June 2024

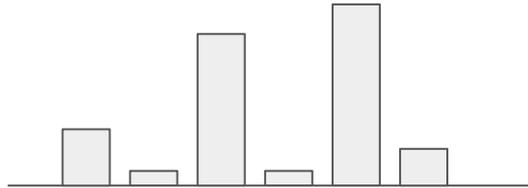
Les grands modèles de langue



*between April 2024 and June 2024

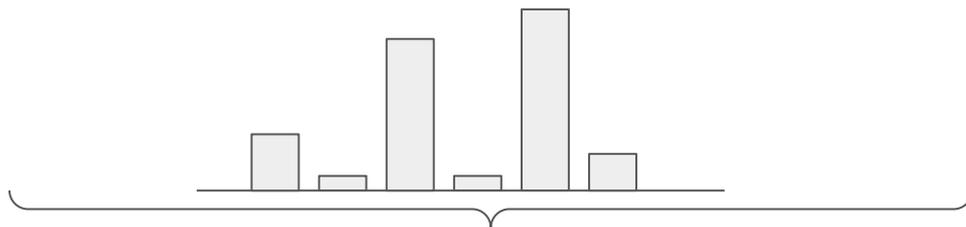
Génération de la langue

*Je mange une carotte .
Je mange une pomme .
Je mange une tarte .
Je mange une pizza .
J'aime bien manger des pommes .
..*



Génération de la langue

Je mange une carotte .
Je mange une pomme .
Je mange une tarte .
Je mange une pizza .
J'aime bien manger des pommes .
..



L'ensemble des **suites de mots possibles** est **très grand** voire **infini!**

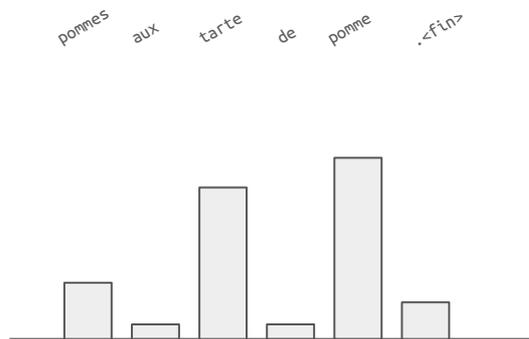
Principe de génération autoregressive

Requête : J'ai mangé une [?]

Principe de génération autoregressive

Requête : **J'ai mangé une [?]**

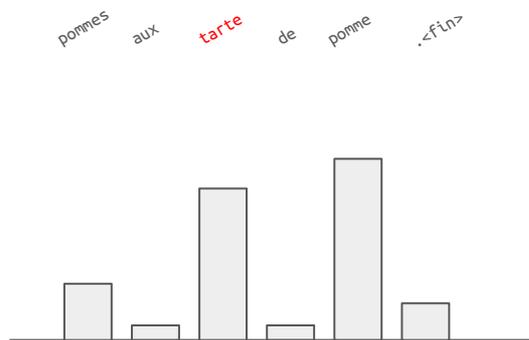
Distribution du prochain mot :



Principe de génération autoregressive

Requête : J'ai mangé une [?]

Distribution du prochain mot :



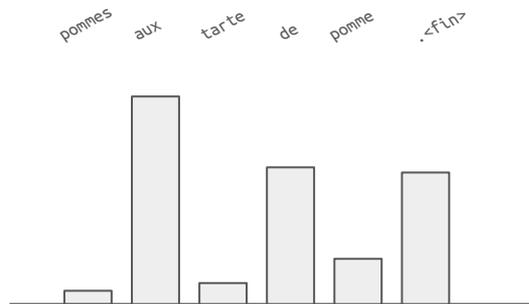
Principe de génération autoregressive

Requête : J'ai mangé une **tarte**

Principe de génération autoregressive

Requête : **J'ai mangé une tarte [?]**

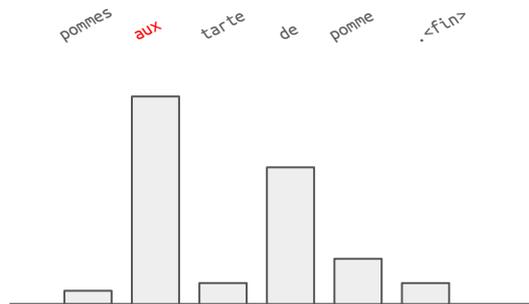
Distribution du prochain mot :



Principe de génération autoregressive

Requête : J'ai mangé une tarte [?]

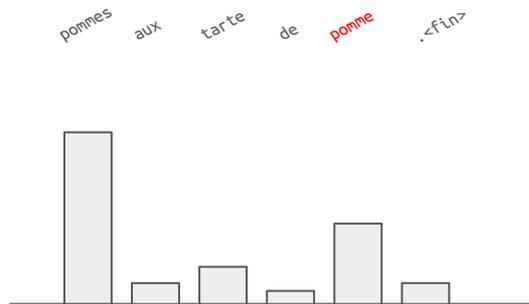
Distribution du prochain mot :



Principe de génération autoregressive

Requête : J'ai mangé une tarte aux [?]

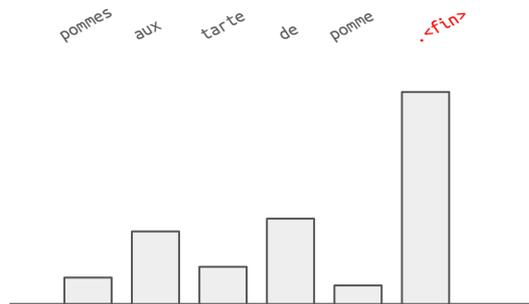
Distribution du prochain mot :



Principe de génération autoregressive

Requête : **J'ai mangé une tarte aux pommes [?]**

Distribution du prochain mot :



Principe de génération autoregressive

Requête initiale : J'ai mangé une

Réponse finale : J'ai mangé une tarte aux pommes .

Principe de génération autoregressive

Fonctionnement :

Pour chaque **préfixe** possible de la langue, on cherche à associer une **distribution** sur le prochain mot à produire.

Ordre de grandeur :

Si le **vocabulaire** contient **30.000 mots** et les préfixes ont une **taille maximale de 10 mots**, le nombre de préfixes possibles est :

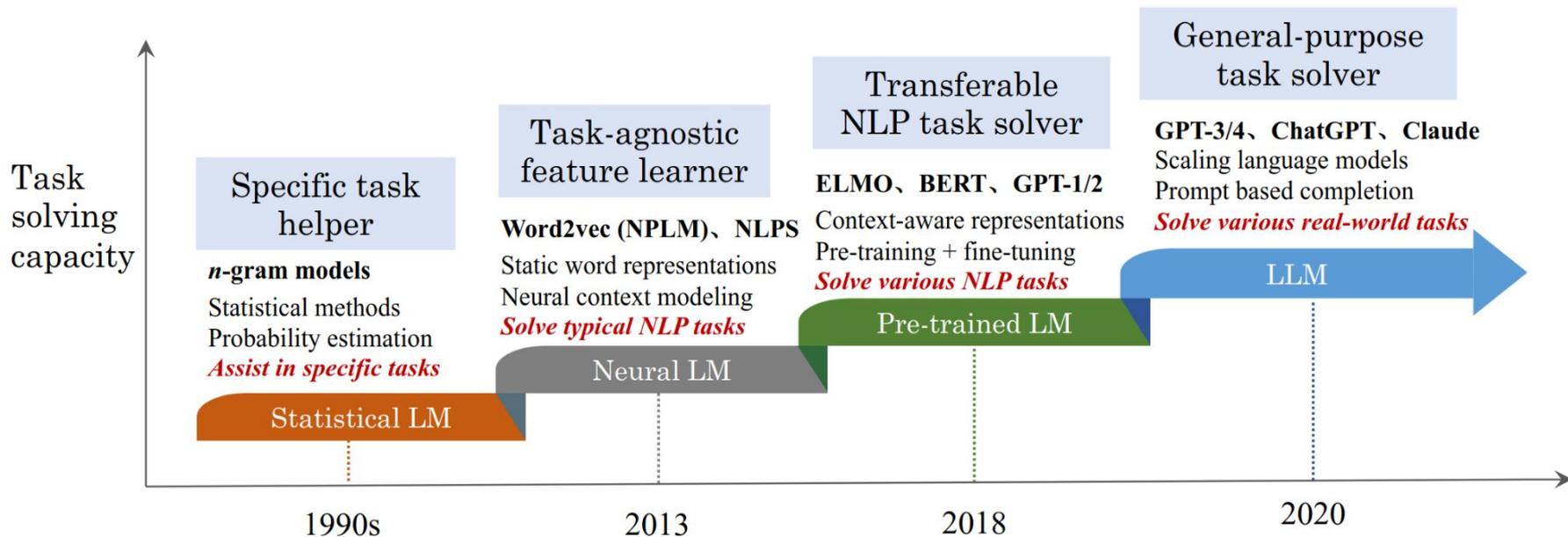
$$30000 + 30000^2 + 30000^3 + \dots + 30000^{10} = 590509683656121870729024300810027000900030000$$

- Beaucoup de préfixes qui n'ont pas de sens réel comme : Je je je je je
- des préfixes peu fréquents dans la langue, etc

Solution :

Apprendre une fonction de distribution conditionnelle approximative -> grand modèle de langue
=> Modèle paramétrique (avec beaucoup de paramètres !)
et non pas une table de proba pour chaque préfixe

Les grands modèles de langue



“Solving” est probablement pas le bon terme à utiliser...

Les grands modèles de langue

Apprentissage :

- Pré-entraînement par **prédiction du prochain mot**.
 - Souvent sur des corpus tirés d'**internet**.
 - Qui contiennent des **biais** et du **langage** qu'on ne veut pas avoir dans la génération.

- Spécialisation sur la tâche d'intérêt e.g. : suivre des instructions.

Requête :

<utilisateur> : Quelle est la capitale de la France ?

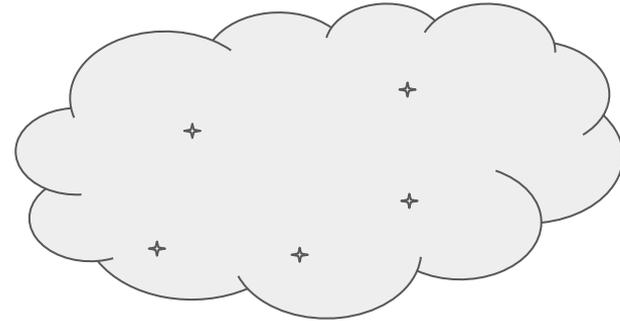
<modèle> : [?]

Réponse :

La capitale de la France est Paris.

Problème de contrôle de la génération

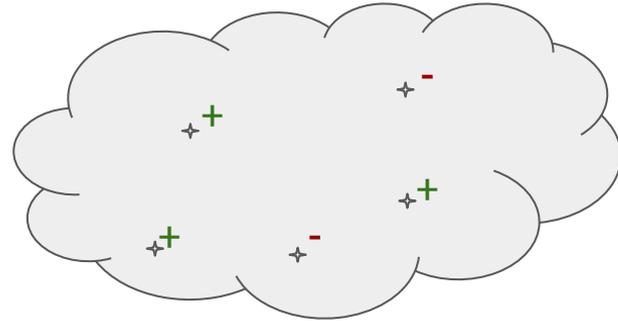
Requête



Plusieurs solutions possibles

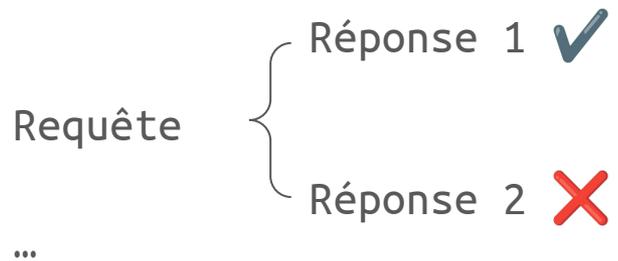
Problème de contrôle de la génération

Requête

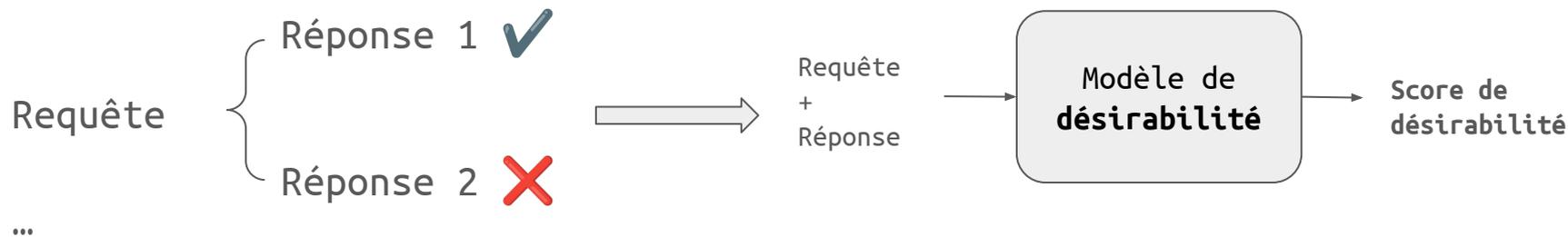


Plusieurs solutions possibles

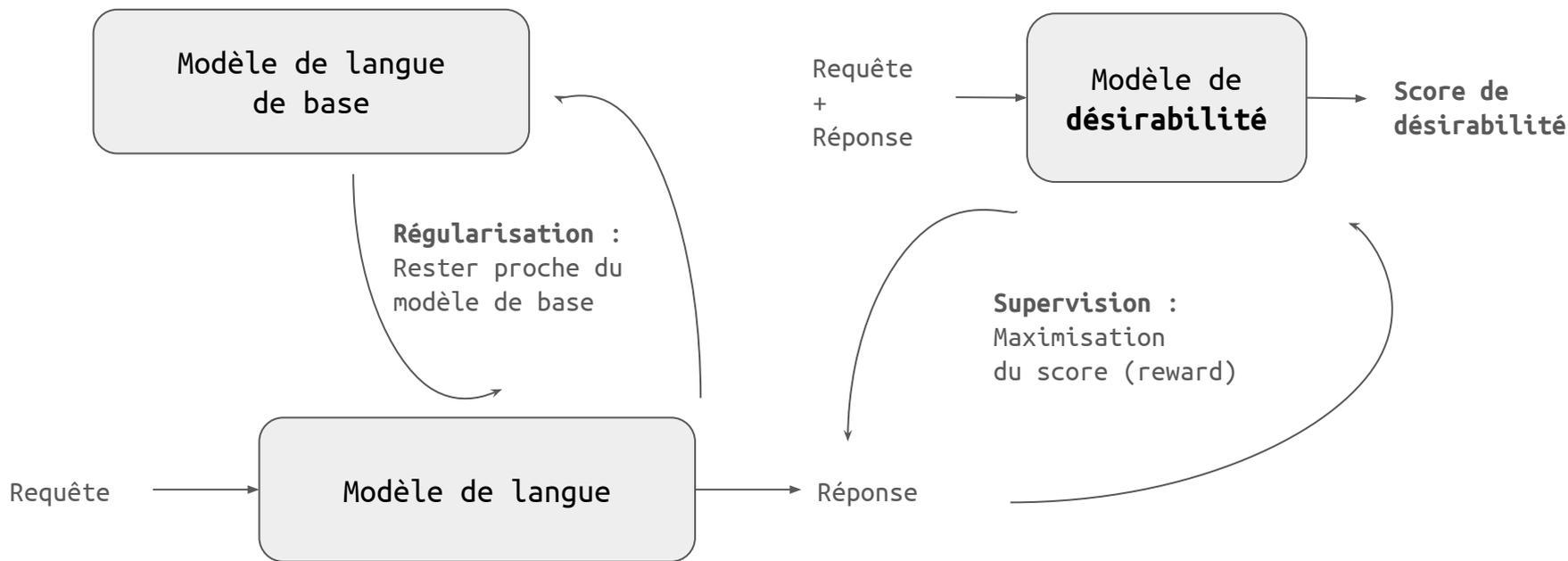
Reinforcement Learning from Human Feedback



Reinforcement Learning from Human Feedback



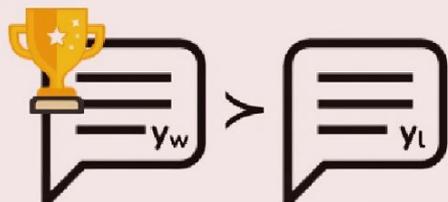
Reinforcement Learning from Human Feedback



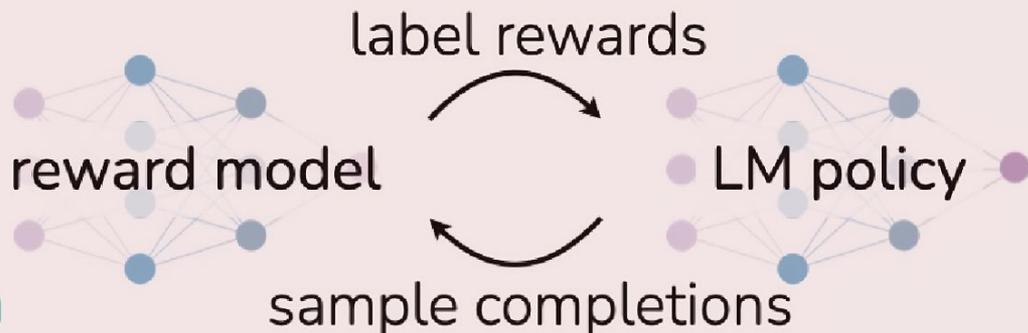
Reinforcement Learning from Human Feedback

Reinforcement Learning from Human Feedback (RLHF)

x: "write me a poem about
the history of jazz"



maximum
likelihood



reinforcement learning

Reinforcement Learning from Human Feedback

Inconvénients:

- Apprentissage par descente de gradient : *SFE trick* pour dériver une estimation de Monte-Carlo du gradient => **apprentissage instable**.
- **Budget de calcul**, plusieurs copies du modèle.
 - => jusqu'à 4 copies !
- **Online** v.s. **offline**.

Objectif similaire aux problèmes d'apprentissage par renforcement :

$$\max_{\pi_{\theta}} \underbrace{\mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_{\theta}(y|x)} [r_{\phi}(x, y)]}_{\text{Maximisation du score de "reward"}} - \underbrace{\beta \mathbb{D}_{\text{KL}} [\pi_{\theta}(y | x) || \pi_{\text{ref}}(y | x)]}_{\text{Minimisation de la distance au modèle de base}}$$

Maximisation du score
de "reward"

Minimisation de la distance
au modèle de base

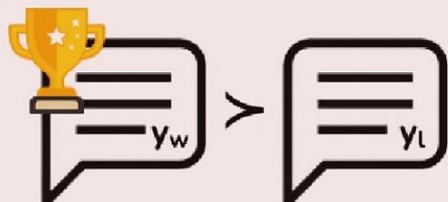
Approximation de Monte-Carlo des espérances:

- Le gradient par rapport aux paramètres du LLM est nul :(
- Le gradient en utilisant le SFE trick est non nul et non biaisé ! :)

Alternative : Direct Preference Optimization

Reinforcement Learning from Human Feedback (RLHF)

x: "write me a poem about
the history of jazz"



preference data

maximum
likelihood



reward model

label rewards



LM policy

sample completions



reinforcement learning

Alternative : Direct Preference Optimization

Direct Preference Optimization (DPO)

x: "write me a poem about
the history of jazz"



maximum
likelihood



Alternative : Direct Preference Optimization

- Récompense implicite, entraînement du modèle de langue directement sur les données de préférences.
- Sans devoir passer par de l'apprentissage par renforcement.

Inconvénients:

- Overfitting des préférences.
- Oubli catastrophique (Catastrophic forgetting).

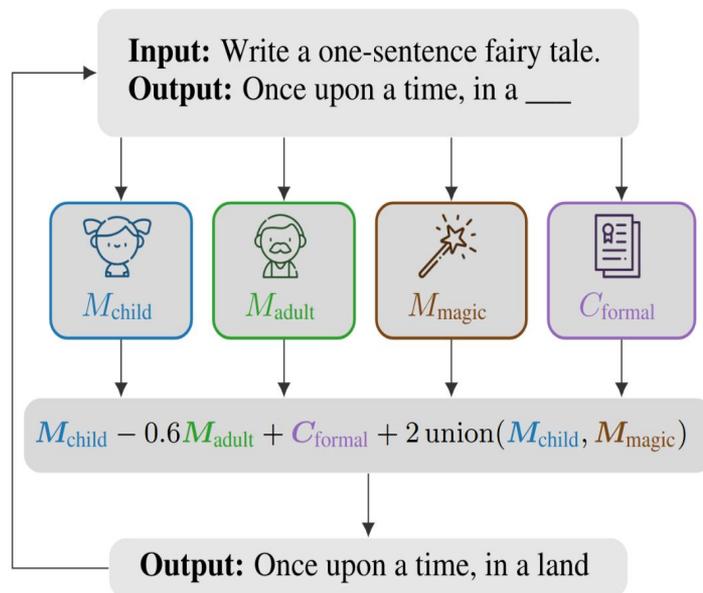
Alternative : Génération contrôlée

E.g. d'approches :

- Arithmétique des modèles
- Reward Augmented Decoding

Coût élevé de génération:

Génération depuis plusieurs modèles au lieu d'un seul.



Plan de travail

Nos pistes de recherche:

1. Une ré-interprétation de l'objectif de RLHF qui permet de dériver une alternative **simple à optimiser**.
2. Un framework théorique solide pour DPO, qui permet de généraliser l'approche à des fonctions de perte qui ne conduisent pas à du catastrophic forgetting.
3. Une approche online guidé par l'incertitude du modèle pour minimiser les coûts d'annotation.

En ce moment:

- Ré-implémentation et expérimentation des méthodes existants
- Mise en place de la piste (1)