

Le principe des modèles de langage n'est **pas nouveau** : ce type d'algorithme est déjà présent quotidiennement avec les **assistants de rédaction des messages** instantanés.

Son fonctionnement consiste à **prédire le mot suivant** à partir d'une instruction (« prompt ») saisie par l'utilisateur.



La **phase d'entraînement** permet d'ajuster les paramètres internes à partir de données textuelles issues du web (non vérifiées).



Une **quantité massive de données et de paramètres** permet d'améliorer les performances de prédiction. On ne peut pas pour autant parler de « compréhension » ou de « personnalité » de l'algorithme, il s'agit d'abord de **réponses produites à partir de régularités statistiques**.



Par ailleurs ces modèles peuvent produire des « **hallucinations** », en inventant de fausses informations, ce qui nécessite la plus grande **vigilance** de la part de l'utilisateur.

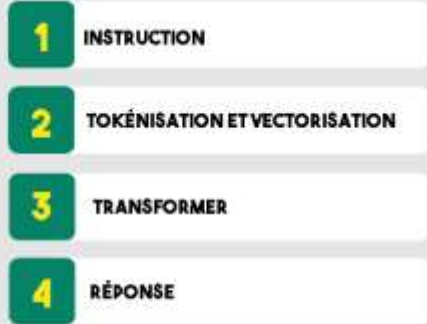


Actuellement les modèles de langage les plus puissants sont détenus par de **grandes entreprises, qui communiquent très peu sur leurs données d'entraînement**.

D'après (Inria Flowers, 2023a)



Le fonctionnement d'un grand modèle de langage (LLM) conversationnel



INSTRUCTION (PROMPT)

L'utilisateur fournit une requête ou une question. Il est recommandé de préciser et contextualiser sa demande pour optimiser la réponse.

PRÉTRAITEMENT ET TOKÉNISATION

La requête est prétraitée pour être comprise par le modèle. Cela inclut la tokenisation : division en unités plus petites, ou tokens (mots, parties de mots, caractères...) et la vectorisation (embeddings).

TRANSFORMER

Les données ainsi encodées sont introduites dans le modèle Transformer, un réseau de neurones de traitement séquentiel du langage avec traitement contextuel (mécanisme d'attention).

RÉPONSE

Le modèle génère une réponse basée sur les régularités statistiques enregistrées dans les masses de données textuelles apprises pendant la phase d'entraînement, et le traitement contextuel de l'instruction utilisateur pour prédire la réponse attendue.

La réponse est décodée en langage naturel pour être rendue compréhensible.



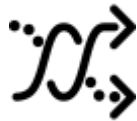
Bilan intermédiaire des tests sur ChatGPT



Réponses aléatoires visant le « vraisemblable »



Pertinentes ou erronées
Parfois « hallucinations »



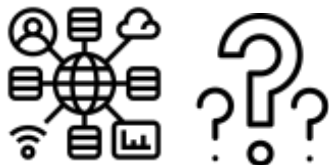
Pas de connexion au web en temps réel



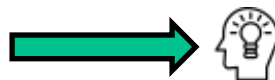
Absence de sources ou sources parfois erronées



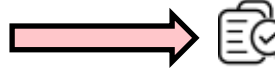
Capacités argumentatives
Réflexives
Parfois de nuances



Sources utilisées ?
Critères à l'origine des réponses ?
Biais culturels, linguistiques...



Potentiel important :
Inspiration – Aide à la rédaction
Traduction - Reformulation
Classification - Synthèse



Grande **vigilance** dans l'utilisation des réponses obtenues