

# IA générative démystifiée

Laure Soulier

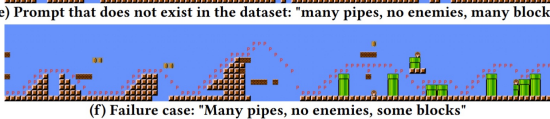
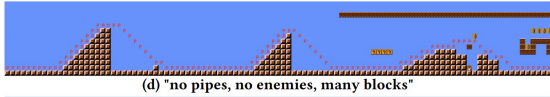
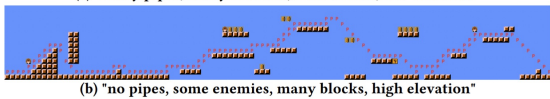
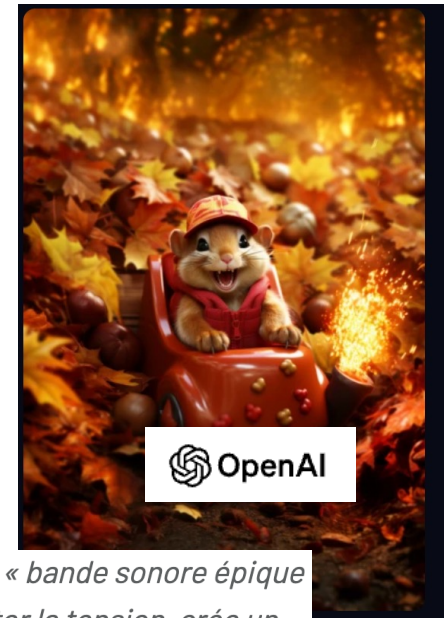
Journée Capsule — 20 juin 2024



Ecris un texte de 5 phrases sur Jian-Yun Nie



Jian-Yun Nie est un chercheur et professeur émérite en informatique. Il est connu pour ses contributions significatives dans le domaine du traitement de l'information et de la recherche d'information. Ses travaux portent notamment sur l'exploration de textes, la recherche d'informations multilingues et la fouille de données. En tant que chercheur de renommée mondiale, Jian-Yun Nie a publié de nombreux articles scientifiques et a participé à des conférences internationales. Ses recherches ont eu un impact important dans le domaine de l'informatique et continuent d'influencer de nouvelles avancées technologiques.



Pour en avoir un aperçu, voici l'extrait généré pour le prompt « *bande sonore épique utilisant des instruments orchestraux. Le morceau fait monter la tension, crée un sentiment d'urgence. Un chœur a cappella chante à l'unisson, ce qui crée un sentiment de puissance et de force* » :

© openAI, Midjourney,  
<https://www.blogdumoderateur.com/musiclm-nouvel-outil-ia-google-generer-musique/>  
And (Sudhakaran et al 2023)



IA, qu'èsaquo?

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**.

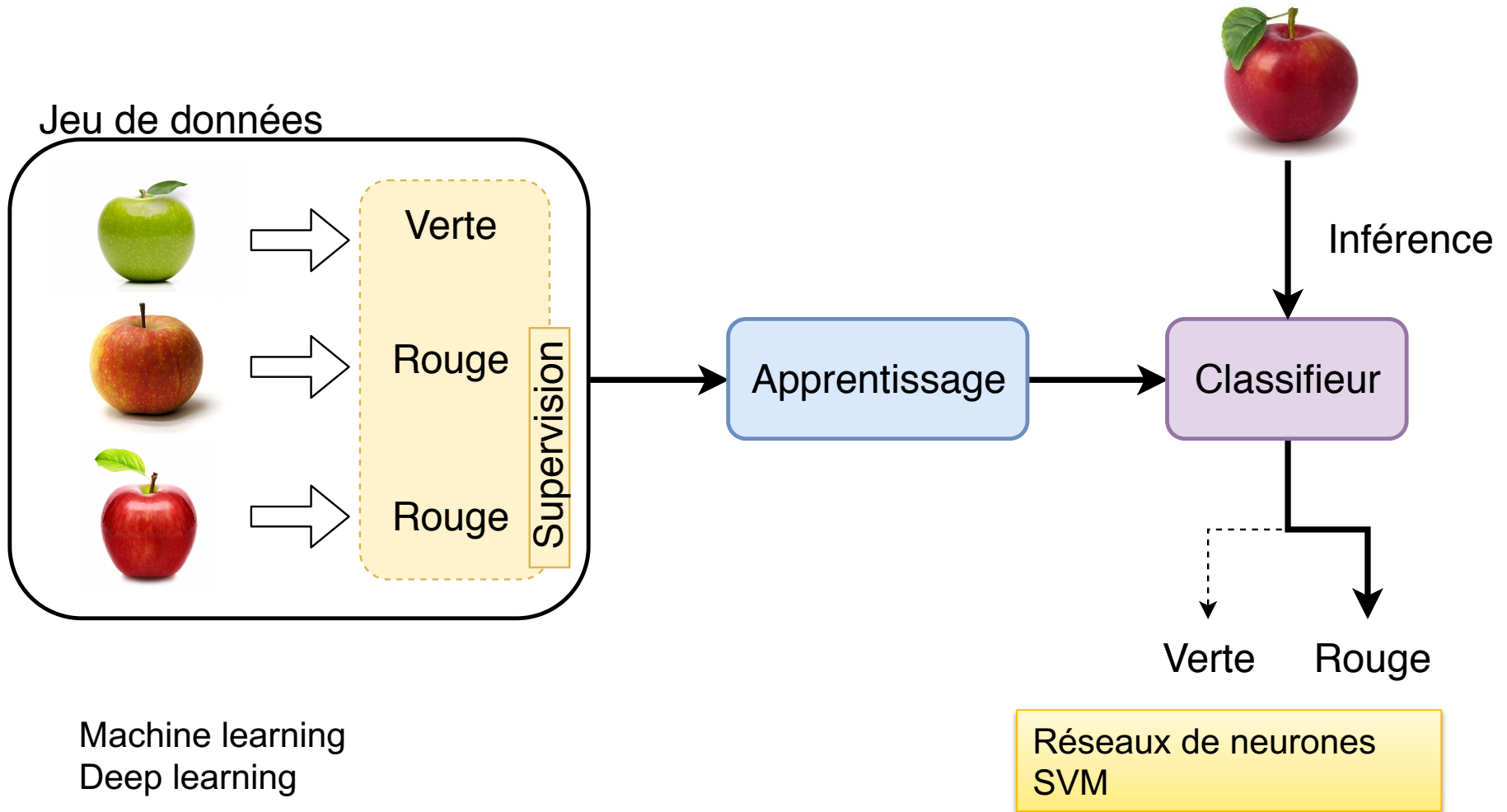
*Marvin Lee Minsky*

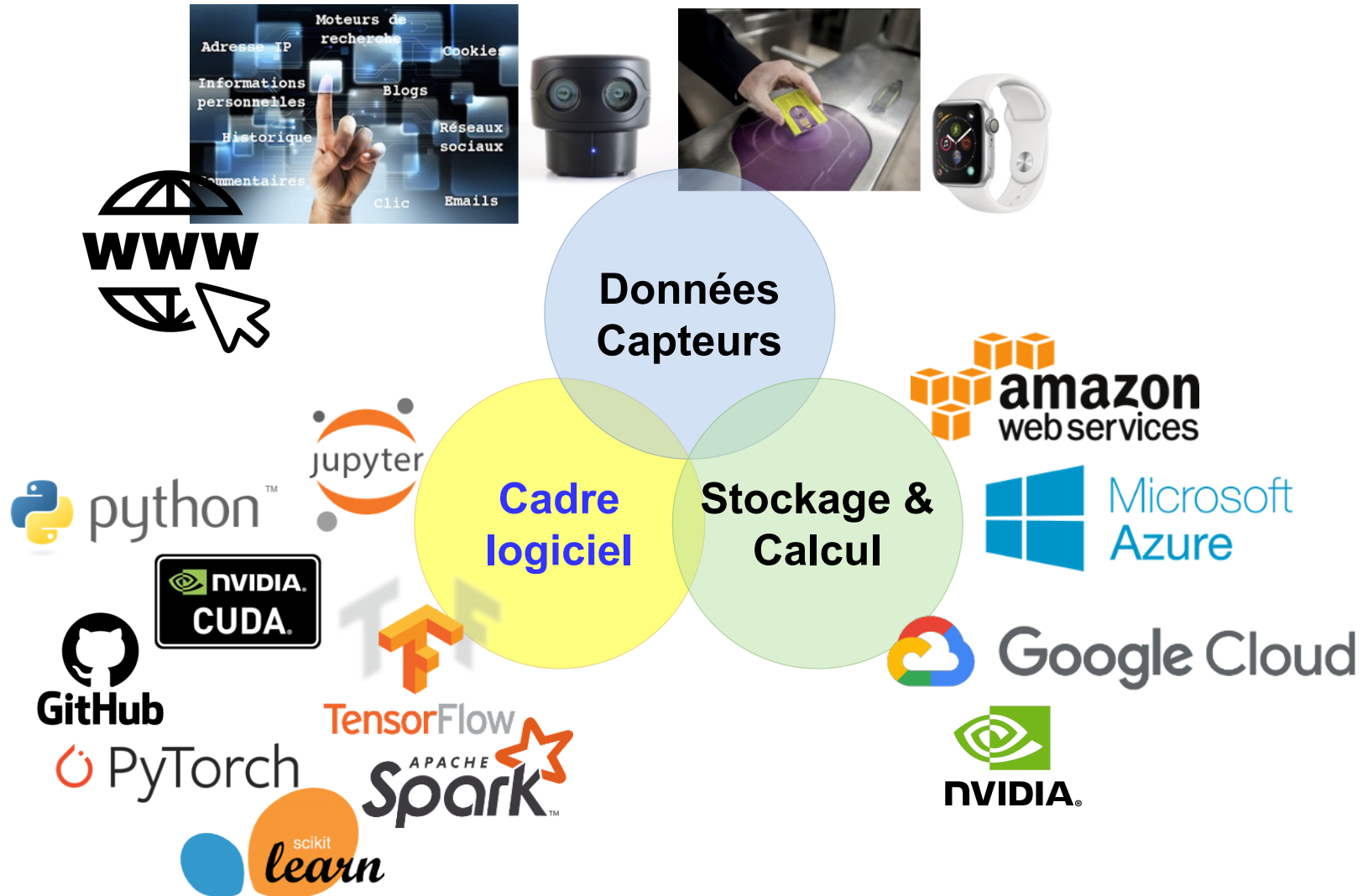
Input (A)	Output (B)	Application
email	spam? (0/1)	spam filtering
audio	text transcript	speech recognition
English	Chinese	machine translation
ad, user info	click? (0/1)	online advertising
image, radar info	position of other cars	self-driving car
image of phone	defect? (0/1)	visual inspection

Ne pas confondre la **NAI (Narrow Artificial Intelligence)**, dédiée à une tâche et la **GAI (General AI)** qui remplace l'humain dans des systèmes complexes.

*Andrew Ng*

→ L'apprentissage par l'exemple (= machine learning)

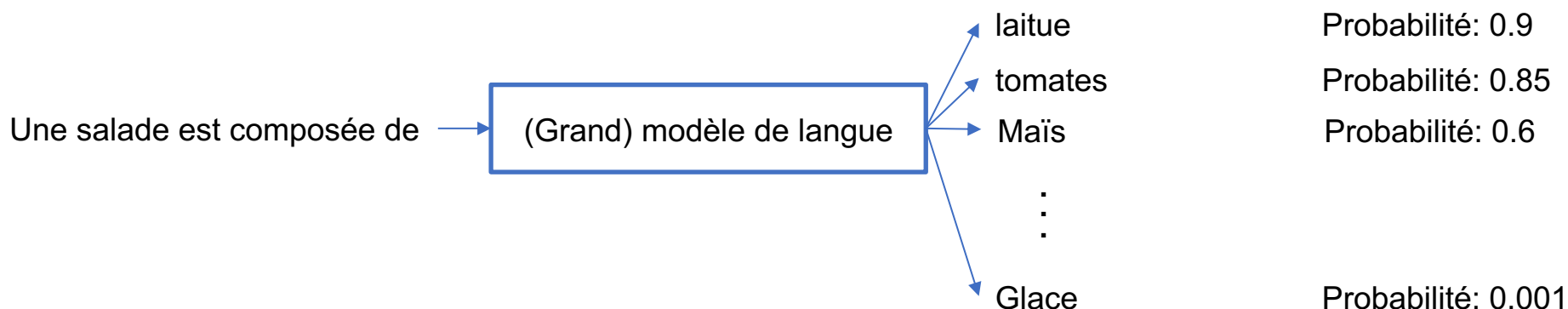




# De la détection d'une pomme à ChatGPT

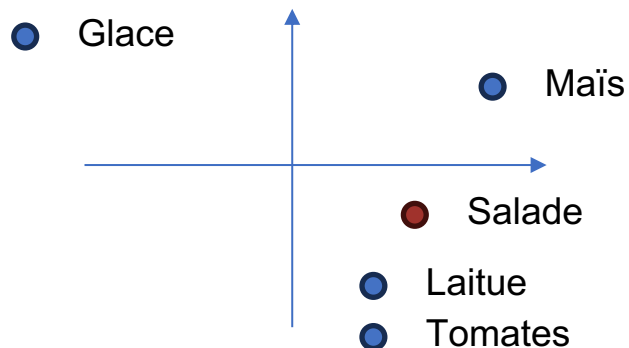
## Principe :

- Modéliser la probabilité d'une séquence  $x_1, x_2, \dots, x_n$
- Les items peuvent être des mots, des caractères, des ngrams/bouts de mots, etc



Etant donné une séquence  $x_1, x_2, \dots, x_{n-1}$ , quelle est la probabilité du prochain item  $x_n$  ?  
 $P(x_n | x_1, x_2, \dots, x_{n-1})$

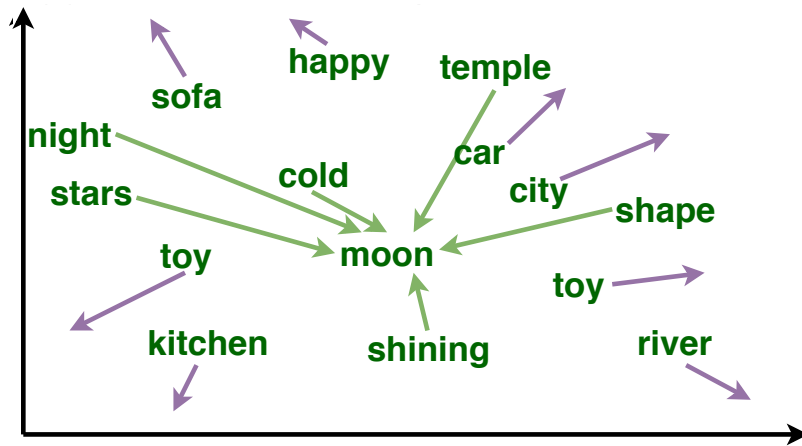
## Sémantique, représentation, espace latent



Salade = (0.3, 0.2, 0.45, -0.1, -0.3)  
Laitue = (0.2, 0.1, 0.38, -0.5, -0.4)  
...  
Glace = (-0.9, -0.3, -0.5, 0.8, 0.7)



## → Algorithme Word2Vec



he curtains open and the moon shining in on the barely  
ars and the cold , close moon " . And neither of the w  
rough the night with the moon shining so brightly , it  
made in the light of the moon . It all boils down , wr  
surely under a crescent moon , thrilled by ice-white  
sun , the seasons of the moon ? Home , alone , Jay pla  
m is dazzling snow , the moon has risen full and cold  
un and the temple of the moon , driving out of the hug  
in the dark and now the moon rises , full and amber a  
bird on the shape of the moon over the trees in front  
But I could n't see the moon or the stars , only the  
rning , with a sliver of moon hanging among the stars  
they love the sun , the moon and the stars . None of  
the light of an enormous moon . The splash of flowing w  
man 's first step on the moon ; various exhibits , aer  
the inevitable piece of moon rock . Housing The Airsh  
oud obscured part of the moon . The Allied guns behind

2000

Modèle pionnier  
de Bengio

2012

Word2Vec,  
FastText, ...

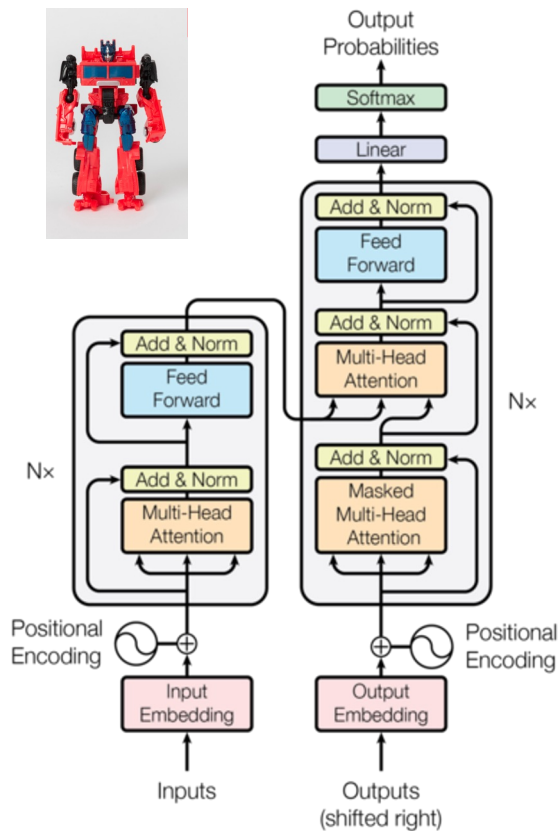
2014

Doc2Vec,  
FastSent, ...

2017

Représentations contextuelles  
Transformer networks  
Bert, T5, GPT, ...

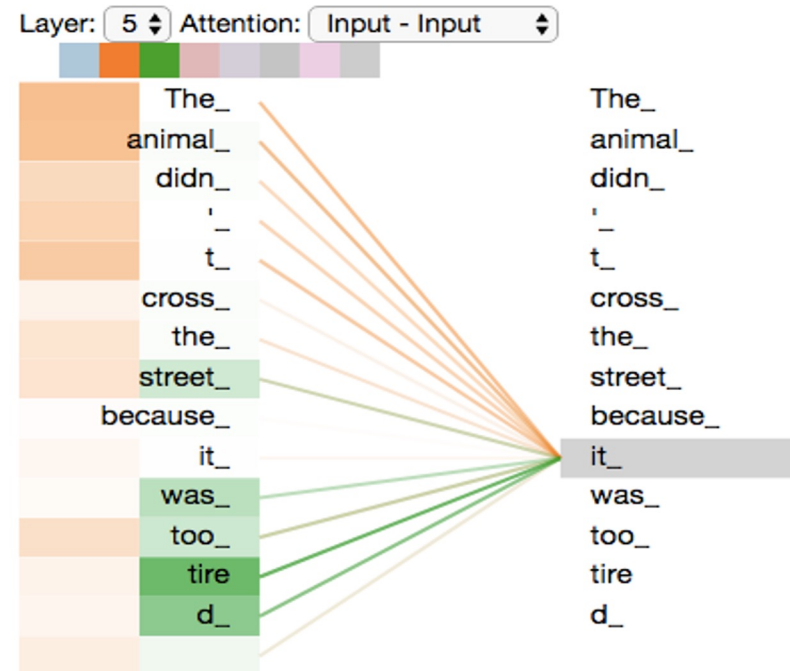
## Transformer (2017)



Un encoder-decoder avec :

- Environ 65 millions de paramètres (maintenant plus)
- Plusieurs blocs successifs
- Des têtes parallèles

... qui estime des représentations contextuelles des items avec l'attention propre

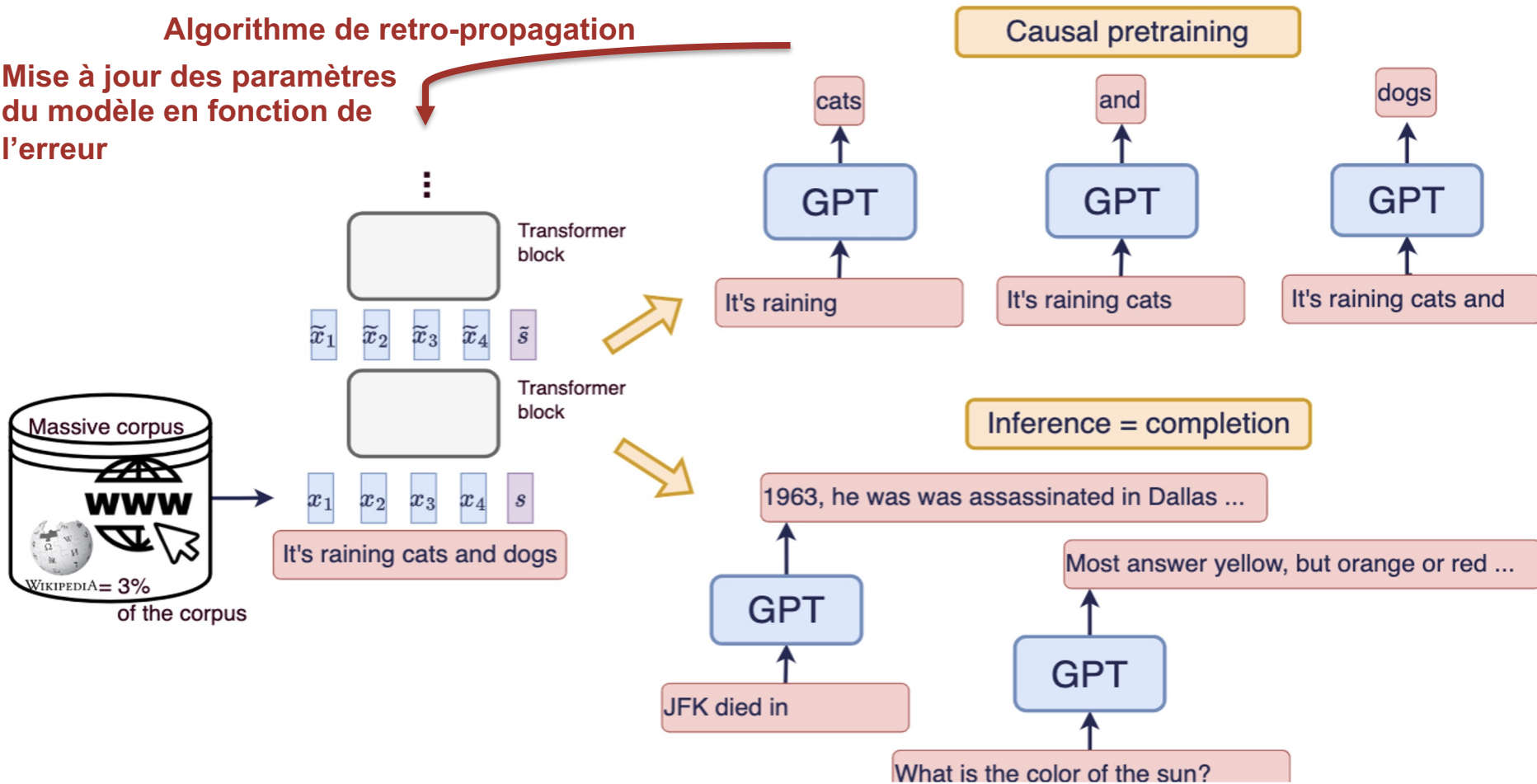


Distinguer *Washington/city* de *Washington/man*

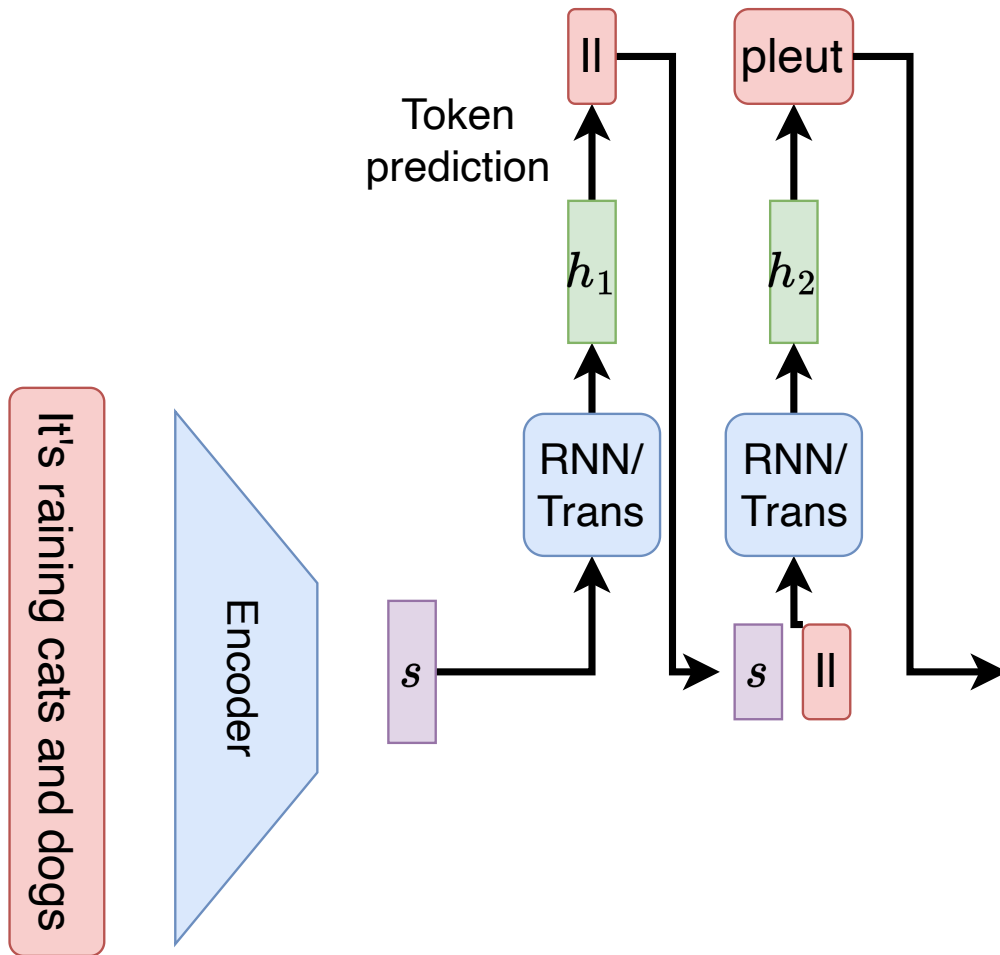
## Entraîner un transformer (e.g. GPT)

### Algorithme de retro-propagation

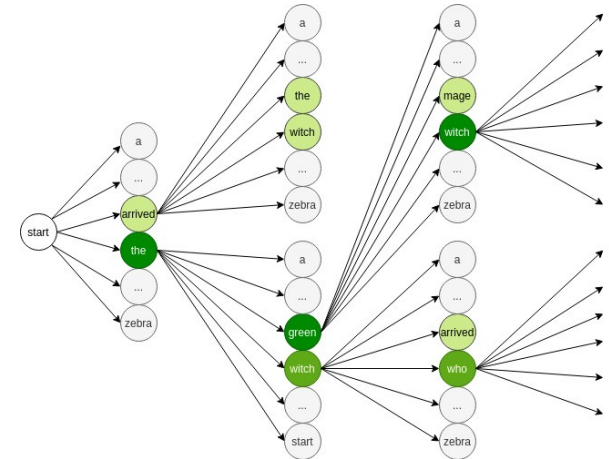
Mise à jour des paramètres du modèle en fonction de l'erreur



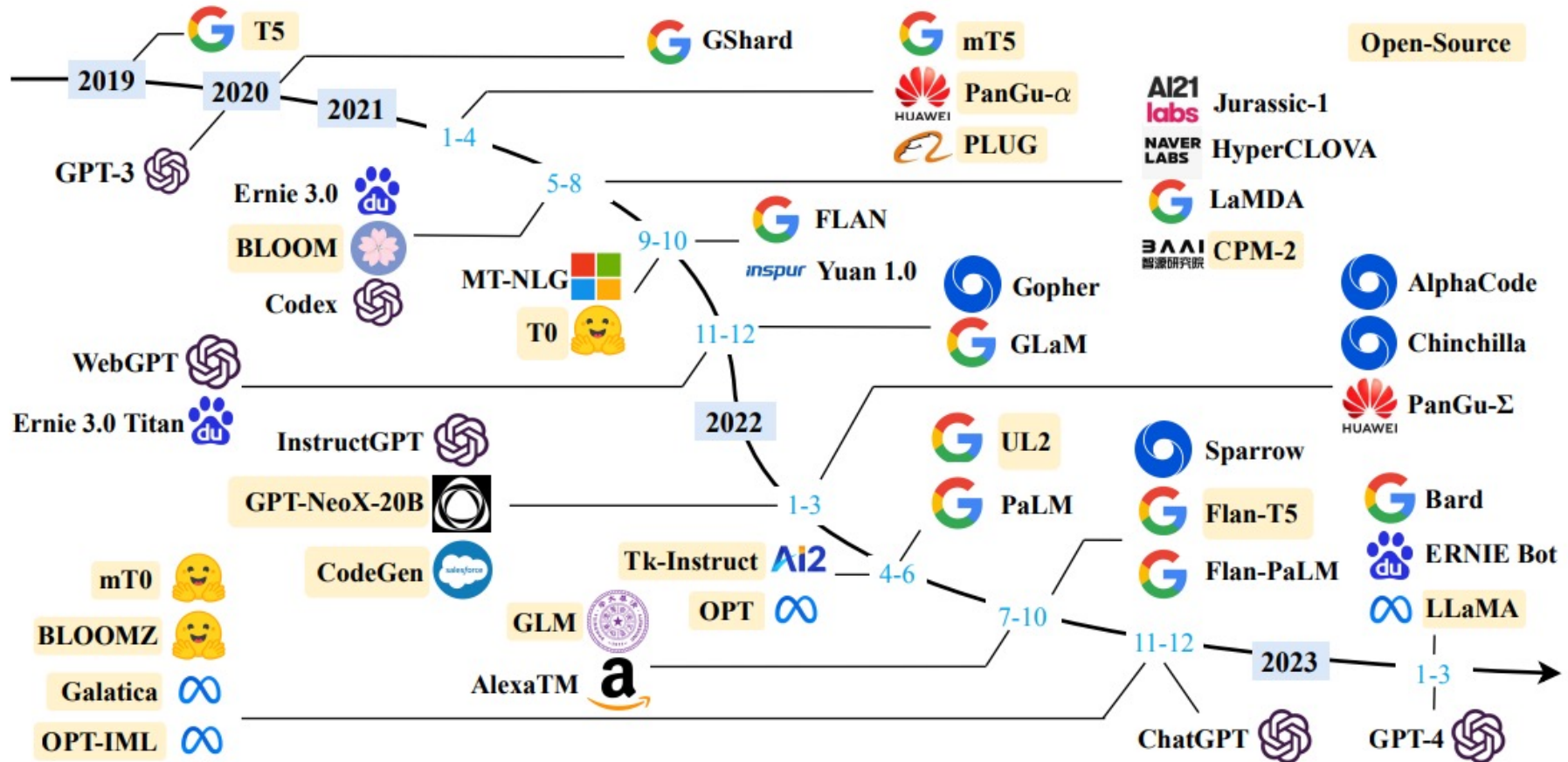
→ Exploiter les architectures précédentes pour écrire du texte



- Génération mot à mot
- Coût très important
- Génération de faisceaux



# De nombreux modèles de langue

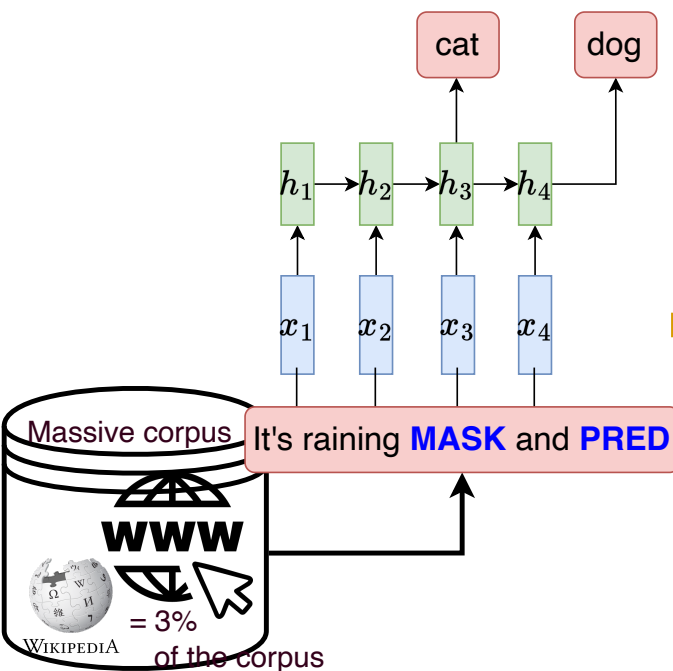


1. Nouvelle tâche
  - Peu de données
  - Choix de la taille des modèles
2. Modèle de langue
  - Connaissances générales
3. Adaptation pour une tâche
  - Traduction
  - Détection d'entités nommées
  - ...

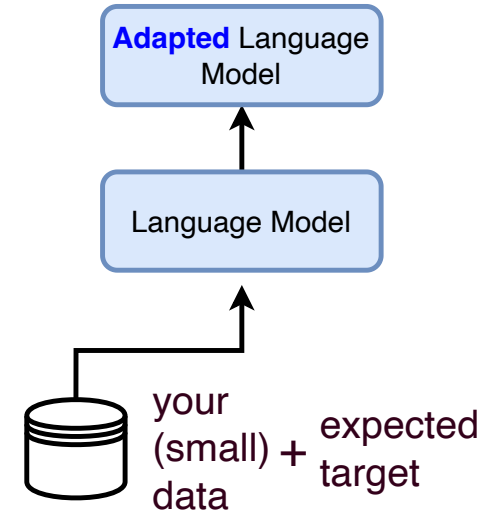
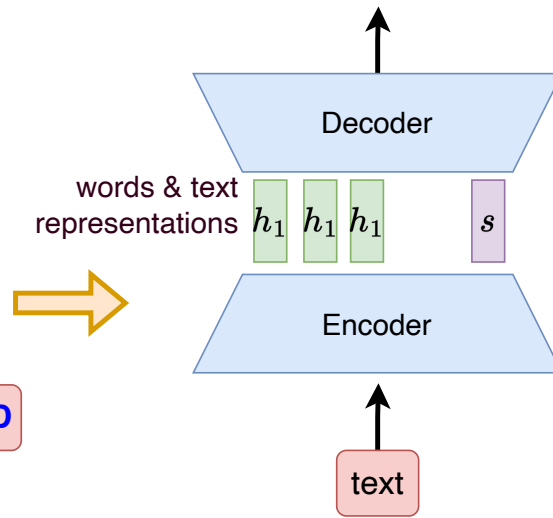
## Pretraining

## Pretrained Language Model

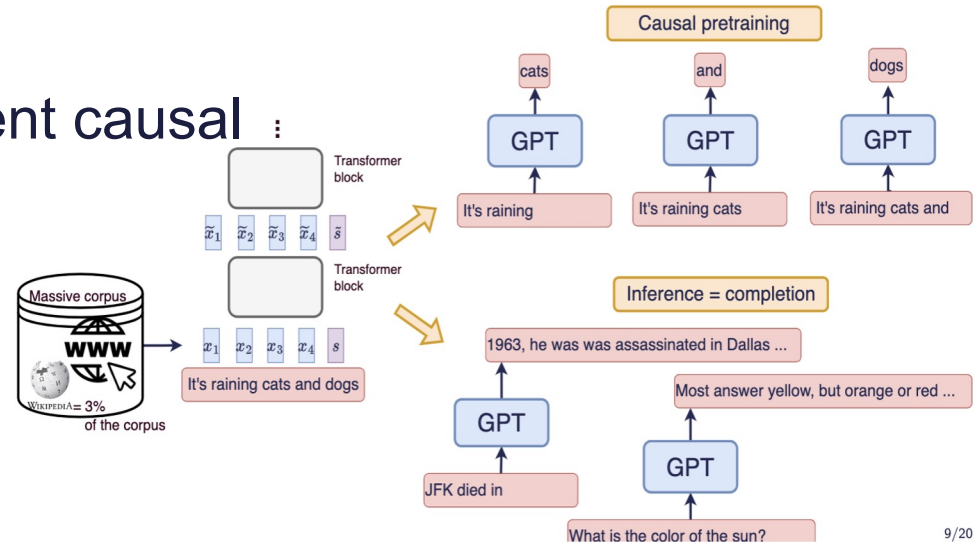
## Finetuned Model



Word prediction; sentence completion; ...



## → Etape 1: Pré-entraînement causal :

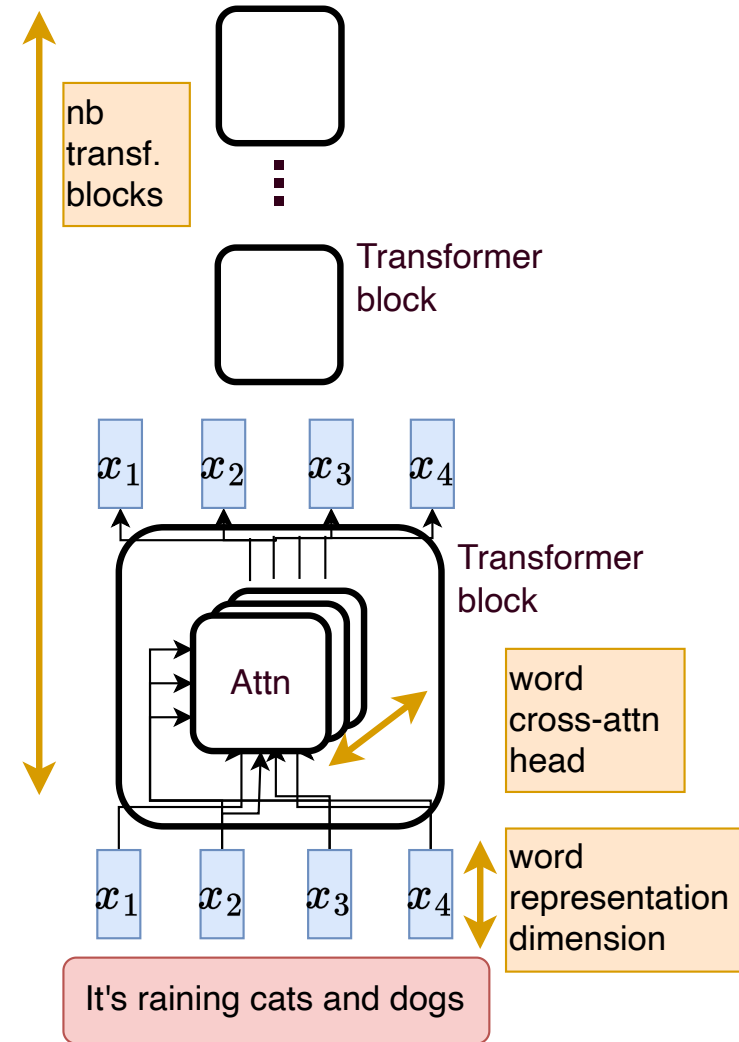


Plus...

- de mots en entrée [500 => 2k, 32k]
- de dimensions pour les mots [500-2k => 12k]
- de têtes d'attention [12 => 96]
- de blocks/couches [5-12 => 96]

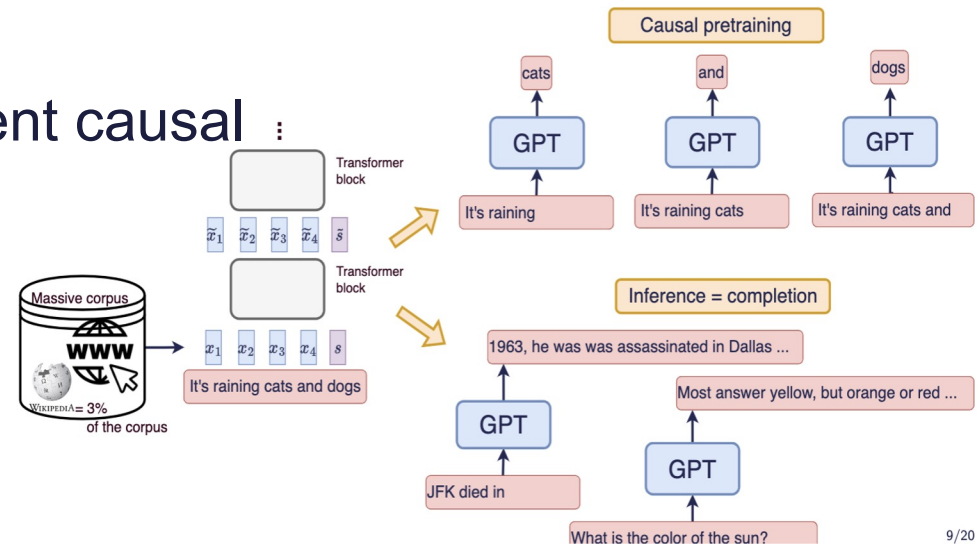
## 175 Milliards de paramètres... Ca fait quoi?

- $1.75 \cdot 10^{11} \Rightarrow 300 \text{ Go} + 100 \text{ Go}$  (stockage des données en inférence)  $\approx 400 \text{ Go}$
- GPU NVidia A100 = 80Go de mémoire (=20k€)
- Coût pour (1) entraînement: 4.6 Millions d'€



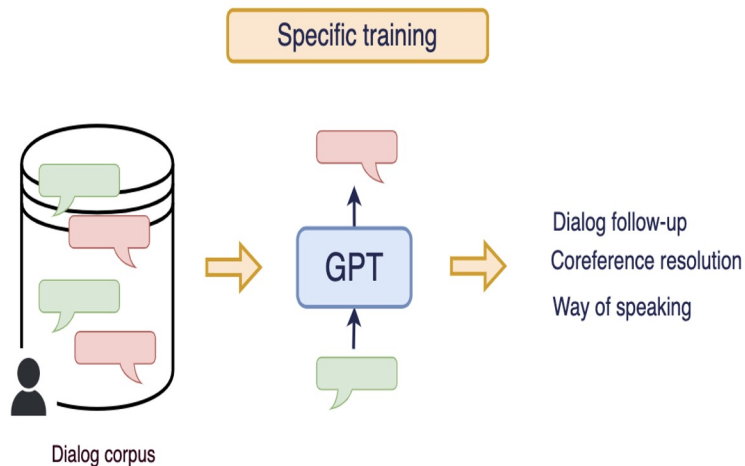


## → Etape 1: Pré-entraînement causal :

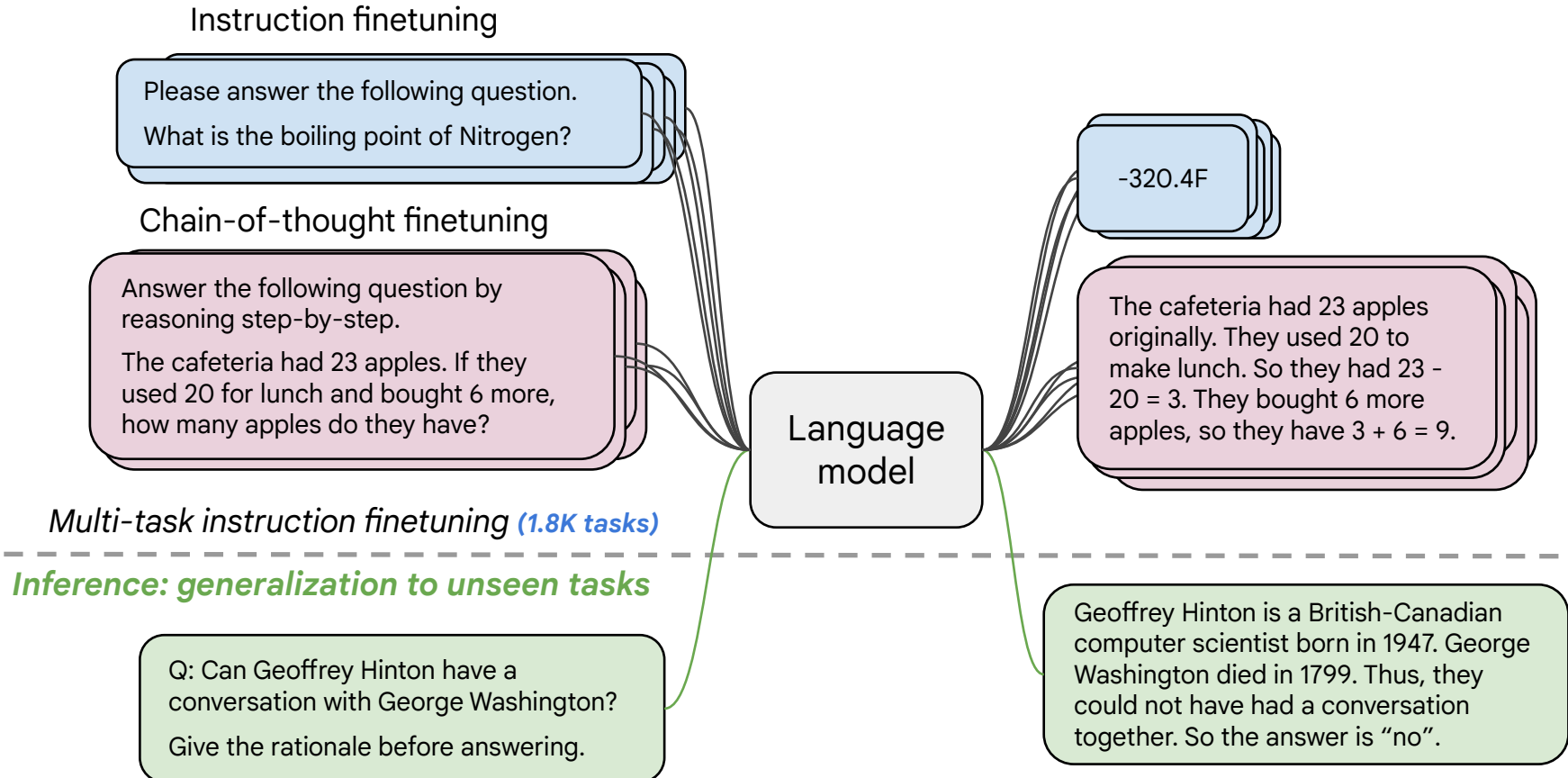


9/20

## → Etape 2: Suivi de dialogue

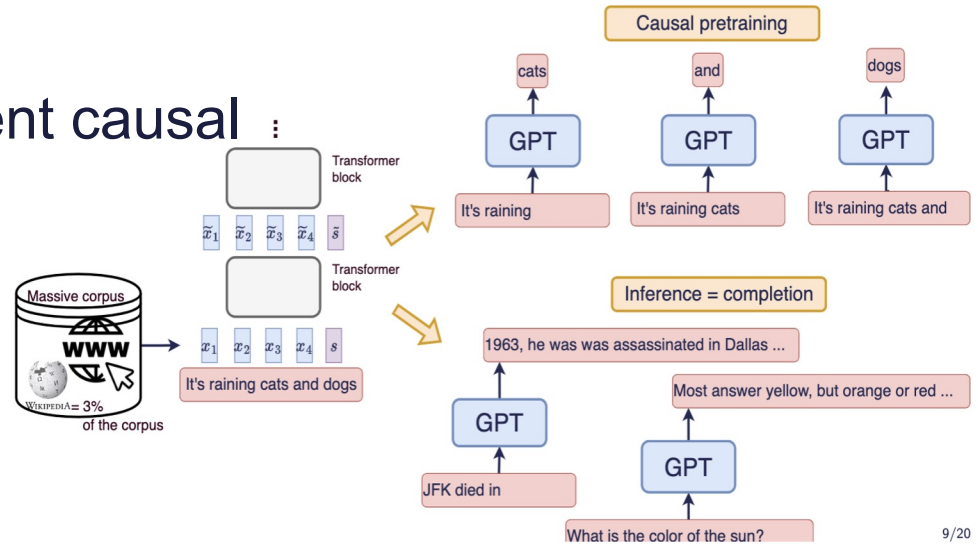


→ Affinage en questions/réponses, raisonnements, ...

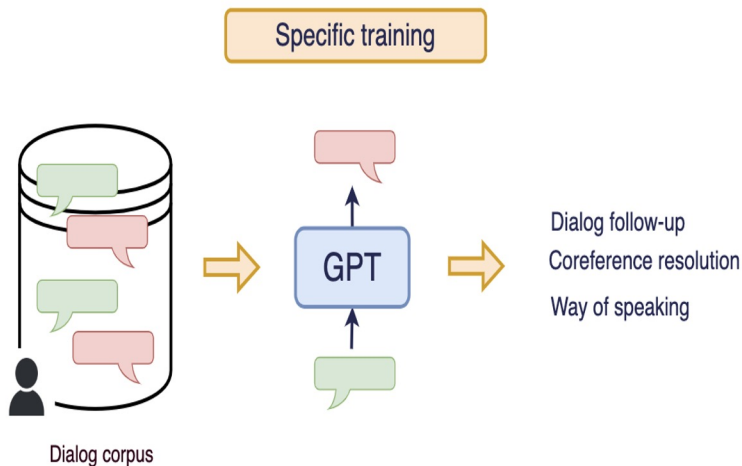


**Importance du prompt**

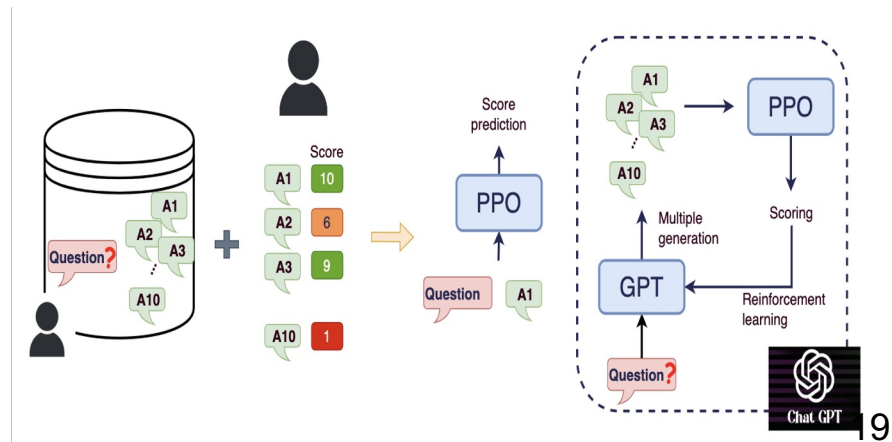
## → Etape 1: Pré-entraînement causal :

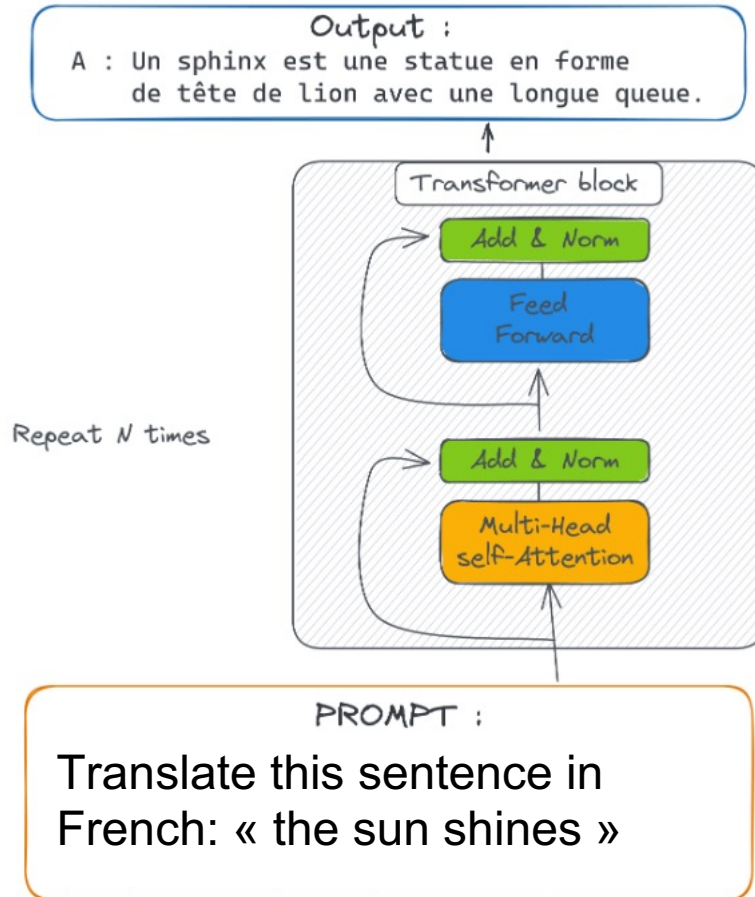


## → Etape 2: Suivi de dialogue



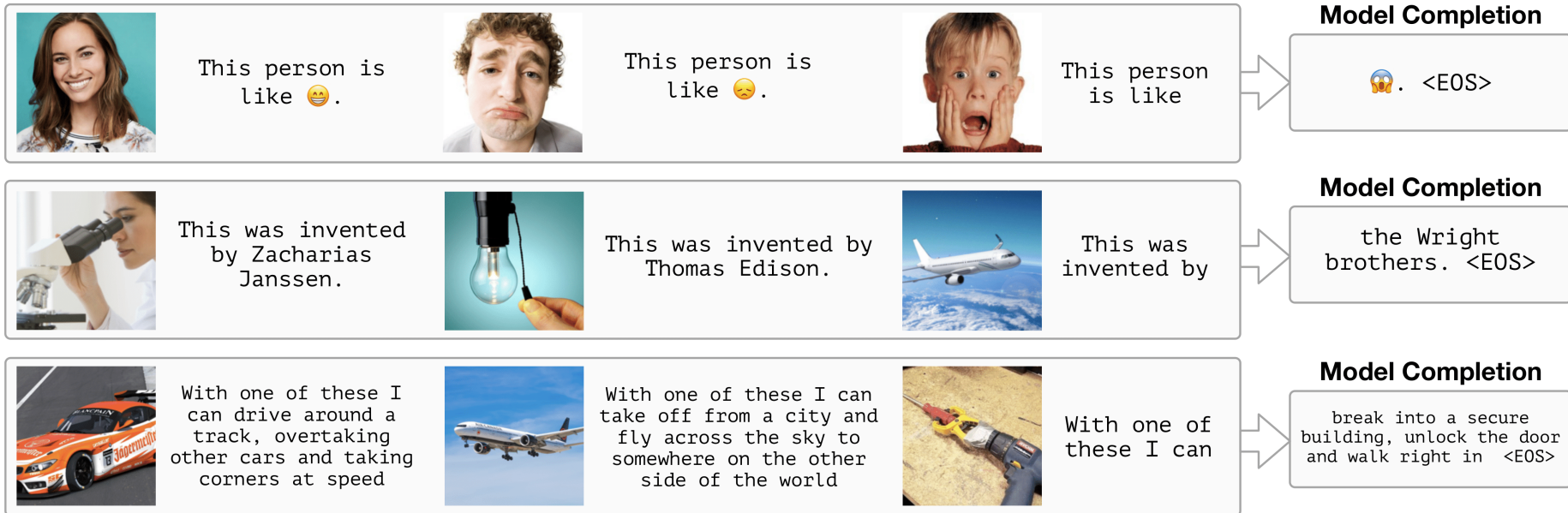
## → Etape 3: Alignement des réponses avec le feedback humain (RLHF)

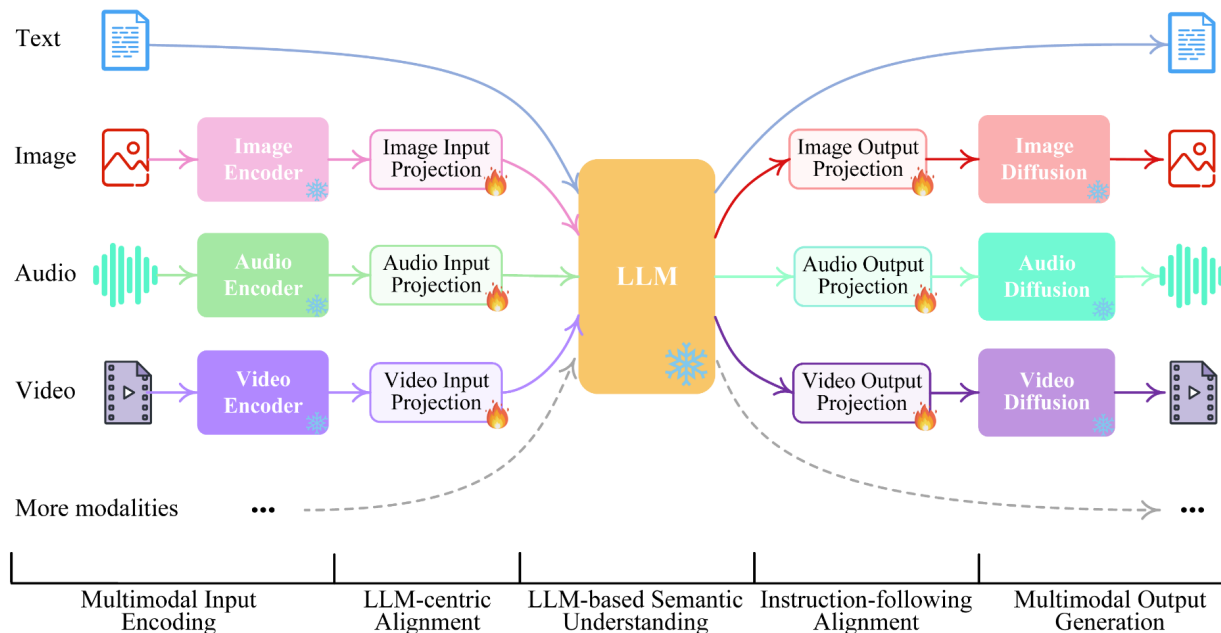
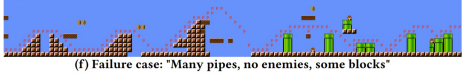
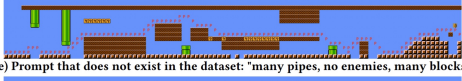
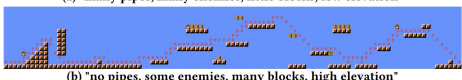
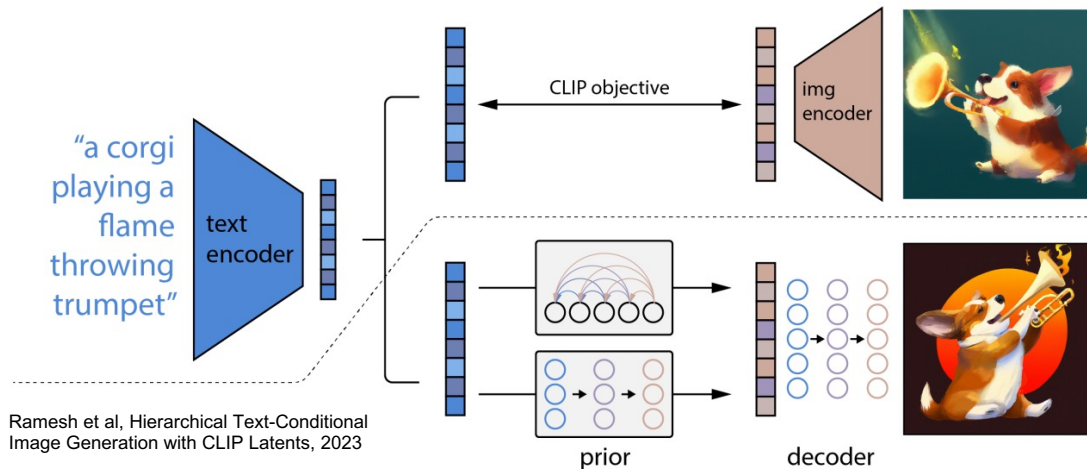




- Prompt : instruction donnée au modèle de langue
- Permet d'explicitier la tâche
- Enjeu : rédiger un prompt précis (contexte, public visé, ...)
- Tout devient génération

- Apprentissage en contexte / In-context learning
- Mentionner des exemples dans le prompt





→ Manque de **véracité / fiabilité**

→ Vérité vs Vraisemblance -  
Génération d'*hallucinations*

→ Incapacité à **s'auto-évaluer**

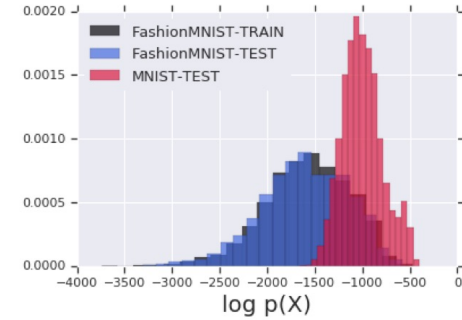
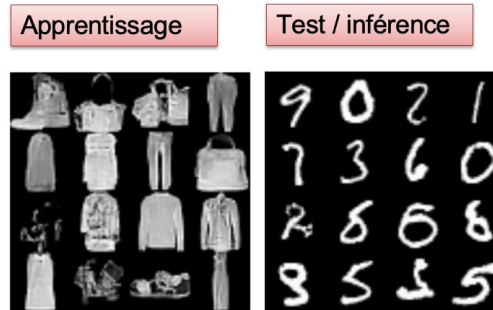
→ **Biais des données**

→ Manque de **stabilité/prédictibilité**

→ *How old is Obama VS how old is Obama?*

→ Manque d'**explicabilité/interprétabilité**

→ Manque de **transparence**



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

Source: 2023 Foundation Model Transparency Index

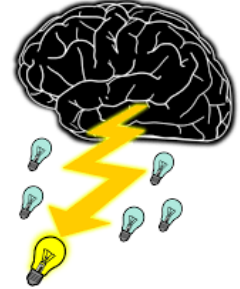
	Meta	Bing	OpenAI	stability.ai	Google	ANTHROPIC	cohere	AI21 labs	Inflection	amazon	Average
	Llama 2	BLOOMZ	GPT-4	Stable Diffusion 2	PaLM 2	Claude 2	Command	Jurassic-2	Inflection-1	Titan Text	
Data	40%	60%	20%	40%	20%	0%	20%	0%	0%	0%	20%
Labor	29%	86%	14%	14%	0%	29%	0%	0%	0%	0%	17%
Compute	57%	14%	14%	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%	33%	30%
Impact	14%	14%	14%	14%	14%	0%	14%	14%	14%	0%	11%
<b>Average</b>	<b>57%</b>	<b>52%</b>	<b>47%</b>	<b>47%</b>	<b>41%</b>	<b>39%</b>	<b>31%</b>	<b>20%</b>	<b>20%</b>	<b>13%</b>	

# Les usages



→ Brainstorming, compte rendu, rédaction de projet

- Développement argumentaire (et recherche de contradiction)
- Mettre en forme des idées
- Reformulation de paragraphes



→ Assistant pour le développement informatique

- Génération de code, recherche d'erreurs, ...



→ Assistant personnel

- Courrier standard, lettres de recommandation, de motivation, de résiliation, ...

→ Assistant pédagogique

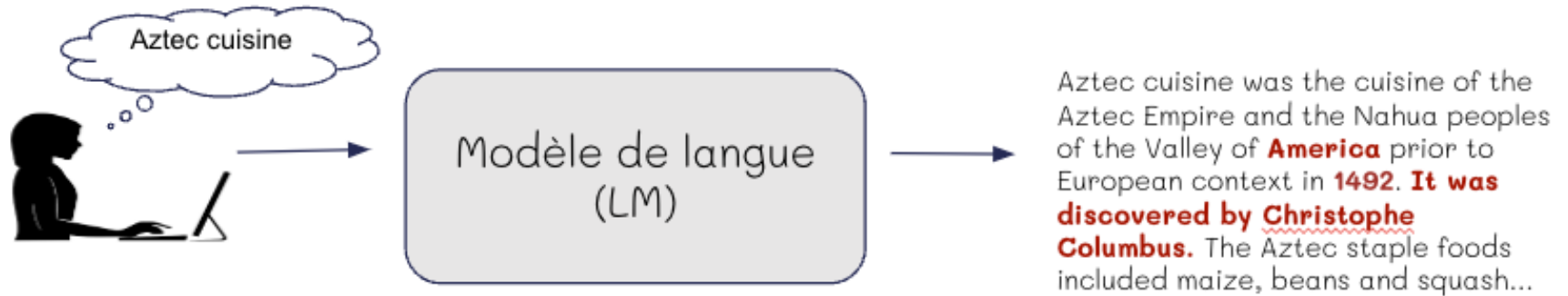
- Wikipédia ++, proposition de plan pour des dissertations, explication de code

→ Analyse de document

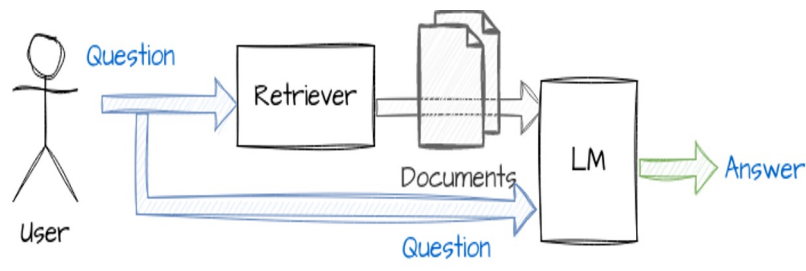
- Extraction d'information, question-réponse, ...



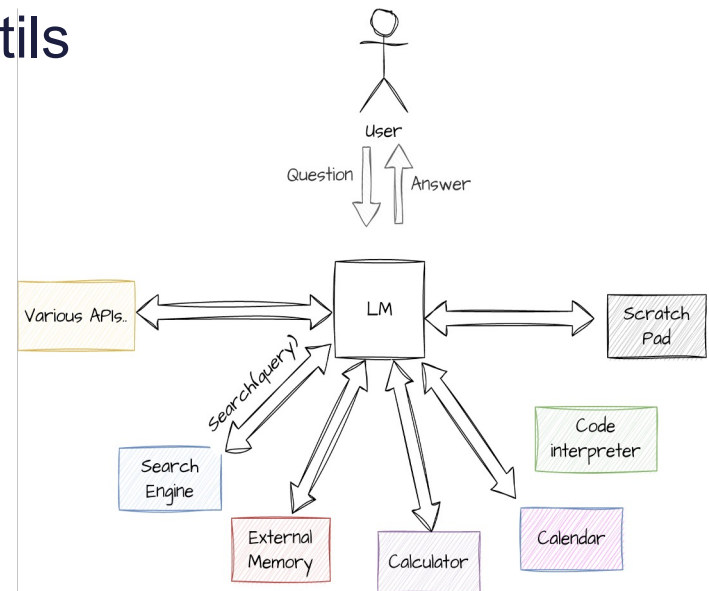
→ Comment limiter les hallucinations ?



→ La génération assistée par les outils



Retrieval-augmented generation (RAG)



- LLM = mémoire (partielle) d'internet
- Champion de la reformulation
- Capacité à comprendre/traduire/générer du code informatique
- Entraîner à répondre à de nombreux types de questions
- Oui, ils vont répondre à beaucoup de choses...  
En faisant régulièrement des petites/grosses fautes



Paradigme de la calculatrice:  
*s'il existe une machine,  
pourquoi apprendre les  
tables de multiplication?*

**Merci !**

---