

Notes - Introduction à l'intelligence artificielle et aux grands modèles de langage

Marius Garénaux-Gruau*, Corentin Val†- IRISA

Table des matières

1	Introduction à la notion de modèle d'IA avec un exemple	2
1.1	Modèles d'apprentissage automatique	3
1.1.1	Consommation électrique en Bretagne	3
1.1.2	Apprentissage Automatique (<i>Machine Learning</i>), entraînement	5
1.2	Quelques exemples de modèles de <i>Machine Learning</i> , et les modèles génératifs	6
1.2.1	IA générative	7
1.2.2	Et les modèles de langage dans tout ça ?	7
1.3	Les impacts des grands modèles génératifs	8
1.3.1	Fabrication des composants et durée de vie du matériel	8
1.3.2	Nettoyage et annotation des données	9
1.3.3	Entraînement des modèles, mise à disposition	9
1.3.4	Impacts sociétaux	9
1.4	<i>Facultatif - Zoom sur les réseaux de neurones</i>	9
2	Explication du fonctionnement des LLM	10
2.1	Modèle de fondation, <i>tokens</i> et machines à camemberts	10
2.2	Entraînement des modèles de fondation	15
2.3	Du modèle de fondation à l'agent conversationnel	17
2.3.1	Ajustement supervisé	18
2.3.2	Ajustement par renforcement	19
3	Comment représenter des <i>tokens</i> avec des nombres ? Les <i>embeddings</i>	20
3.1	Vecteurs	20
3.2	Embedding	21
4	Écosystème et historique des LLM	22
4.1	Différencier les LLM	22
4.2	Historique des LLM	23
4.2.1	1990-2010 : Approches statistiques et neuronales	23
4.2.2	2011-2016 : Word Embeddings et avancées RNN	25
4.2.3	2017-2019 : Apparition des Transformers	26
4.2.4	2020-2021 : Montée en puissance et adoption	26
4.2.5	2022 : Démocratisation et <i>open-source</i>	27
4.2.6	2023-2024 : Multimodalité et spécialisation	27
4.3	Accès et déploiement des LLM	28

*maris.garenaux-gruau@irisa.fr

†corentin.val@irisa.fr



A whimsical and creative image depicting a hybrid creature that is a mix of a waffle and a hippopotamus. This imaginative creature features the distinctive, bulky body of a hippo, but with a texture and appearance resembling a golden-brown, crispy waffle. The creature might have elements like waffle squares across its skin and a syrup-like sheen. It's set in a surreal environment that playfully combines a natural water habitat of a hippo with elements of a breakfast table setting, possibly including oversized utensils or plates in the background. The image should evoke a sense of playful absurdity and culinary fantasy.

Figure 1 – Un exemple d'image générée à partir d'un texte par Stable Diffusion (1)

*Ces notes viennent en complément d'une présentation donnée lors d'une formation d'une journée sur les LLM. Ce document disponible sous la licence **CC-BY-NC-SA**.*

Voici un rapide condensé de ressources à consulter; de la plus abordable à la plus approfondie. Elles ont été utilisées, entre autres, pour rédiger ce document. Vous pouvez les utiliser comme des compléments à ces notes, pour approfondir certains sujets ou tout simplement voir différents points de vue.

- (2) : Science Etonnante ; Vidéo accessible à un large public et complète sur ChatGPT.
- (3) : Rapport du Sénat sur l'IA; très complet, abordable et récent.
- (4) : Andrej Karpathy, présentation abordable et complète (en Anglais) du fonctionnement des LLM.
- (5) : Formation FIDLE ; en live ou rediffusion. Couvre un large éventail de notions, avec les TP correspondants.
- (6) : Panoram'IA : Emission mensuelle de l'IDRIS sur l'actualité en IA
- (7) : Umar Jamil, beaucoup de ressources vidéo, et des exemples de codes sur le fonctionnement détaillé de nombreux modèles. Niveau relativement avancé, mais vaut le détour.
- (8) : Article de recherche récent couvrant assez largement les LLM.

1 Introduction à la notion de modèle d'IA avec un exemple

La notion d'Intelligence Artificielle est de plus en plus présente dans le quotidien de nombreuses personnes. Des algorithmes génératifs (génération d'image à partir de texte, agents conversationnels, ...) aux systèmes de recommandation (Spotify, Instagram, ...); les exemples sont nombreux et de plus en plus impressionnants. Afin de prendre les décisions les plus appropriées (comment utiliser - ou **ne pas utiliser** - ces modèles, lesquels choisir dans quel cadre, ...), il est nécessaire de bien comprendre les **enjeux** et les **impacts** liés à l'utilisation et au déploiement de ces outils.

Il est difficile de trouver une définition rigoureuse et précise de l'Intelligence Artificielle. Pour cause, on peut aborder cette notion sous de nombreux aspects : scientifique, économique, culturel, technologique, juridique,... Trouver une définition englobant tous ces exemples n'est alors pas si simple... On se réfère au très bon rapport de l'OPECST (Office Parlementaire d'Evaluation des Choix Scientifiques et Technologiques)

sur les nouveaux développements de l'intelligence artificielle (3), page 41 (1.B.1); pour une explication plus détaillée sur ce sujet.

On gardera à l'esprit dans ce document que les modèles d'Intelligence Artificielle dont on va parler sont d'abord des **programmes informatiques**. Ils permettent d'automatiser des tâches.

1.1 Modèles d'apprentissage automatique

Nous préférons dans la suite restreindre la discussion aux **modèles d'Apprentissage Automatique** (*Machine Learning*). Ces modèles ont une définition plus précise, que l'on donne juste après un exemple : comment modéliser le lien entre consommation électrique et température extérieure ?

1.1.1 Consommation électrique en Bretagne

Nous avons téléchargé des données de consommation quotidienne (9) et de température moyenne (10) en Bretagne, depuis 2016. Ces données peuvent être affichées dans un tableau à 3 colonnes : la date, la température et la consommation électrique. Ce tableau peut être représenté sur un graphique, en affichant un point pour chaque relevé. La coordonnée sur l'axe horizontal est la température du jour, et celle sur l'axe vertical est la consommation totale quotidienne en Bretagne. Ceci produit un nuage de points illustrant le lien entre consommation électrique et température (Figure 2).

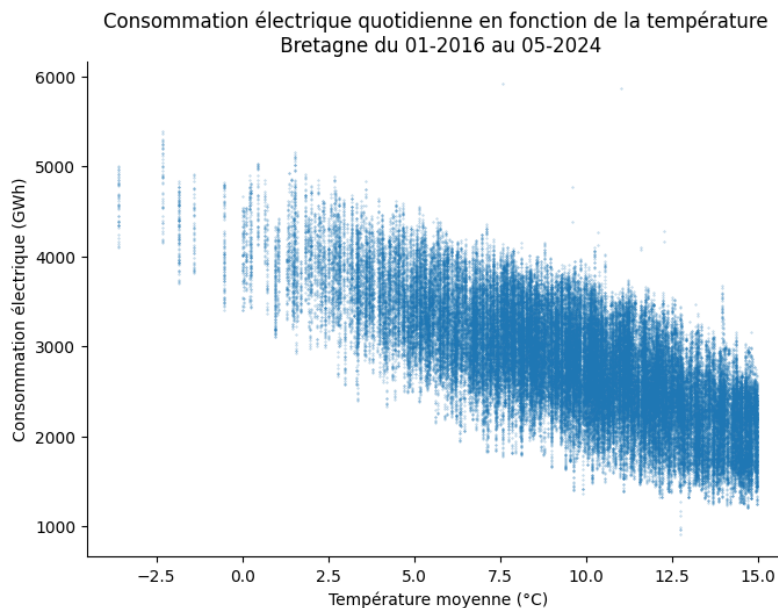


Figure 2 – Consommation électrique et température en Bretagne depuis 2016

Imaginons désormais que l'on souhaite avoir une idée de la consommation électrique pour toute la Bretagne demain ; en sachant qu'il fera en moyenne 10°C (avec une prédiction météo par exemple). En utilisant les données récoltées sur les 10 dernières années, on peut construire un programme informatique qui nous permet de faire cette prédiction.

Exemples de modélisations

Il y a plusieurs manières de réutiliser ces données pour cette tâche de prédiction. Ces différentes options sont **différentes modélisations** du lien entre température et consommation électrique en Bretagne. En voici une :

Régression linéaire

On observe une tendance globale : lorsque la **température diminue**, la **consommation augmente** à un rythme **constant**. C'est très probablement lié au chauffage - mais on ne cherche pas spécialement à comprendre ce lien ici. Forts de cette observation, on **choisit de modéliser** le lien entre température et consommation par une **droite**. Une **droite**, oui, mais laquelle ? On peut en tracer une infinité, en changeant la **hauteur de la droite**, et sa **pen**te.

Petit rappel : une droite a pour équation $y = a \cdot x + b$. Cette équation décrit le lien entre abscisse (x) et ordonnée (y), avec **deux paramètres** : **a** et **b** - la pente et la hauteur de la droite. Faire varier ces deux paramètres change la droite.

Parmi toutes les droites possibles, certaines sont *plus proches* du nuage de point que d'autres. Et donc, certaines droites vont mieux modéliser le phénomène voulu que d'autres. Pour avoir un modèle pertinent, il faut définir cette notion de "**proche**"¹, et **trouver** des paramètres a et b qui nous fournissent une droite la plus proche possible du nuage de points. Voir Figure 3 la droite obtenue.

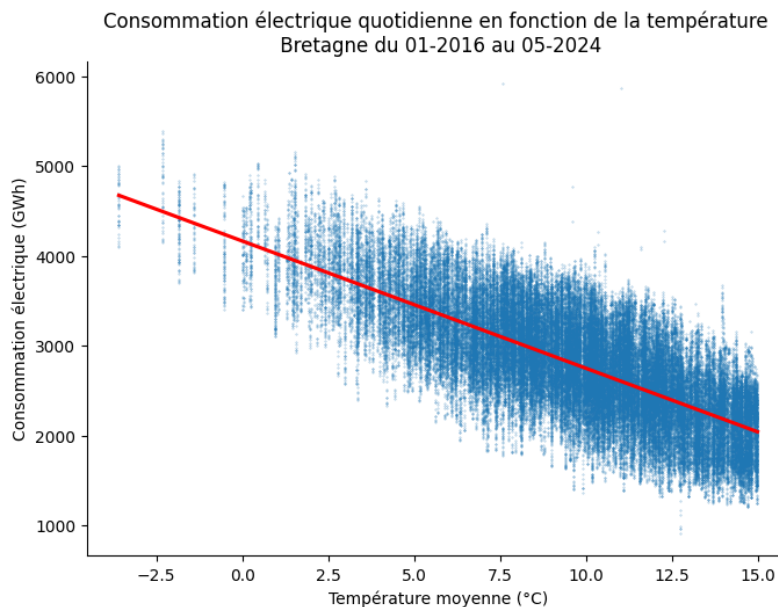


Figure 3 – Consommation électrique et température en Bretagne depuis 2016

Selon ce modèle, s'il fait 10°C demain, la consommation moyenne sera d'environ 2800 GWh. Avec le modèle, on a une réponse à notre question. Voici un second exemple de choix de modèle :

Plus proches voisins

Sans se soucier de la tendance globale, on peut aussi avoir l'idée de faire la prédiction en regardant directement dans les données. On peut ainsi construire un programme informatique qui, étant donnée une température (par exemple 10°C), calcule la moyenne des consommations électriques sur les 10 éléments du jeu de données dont les températures étaient les plus proches de 10°C. C'est aussi un modèle ! Il se base uniquement sur les données. Il a également des paramètres que l'on peut modifier (le nombre de

1. Bien qu'elle semble floue, la notion de "proche du nuage de point" a un sens très précis ici. Pour un point donné, on peut calculer sa différence de hauteur avec la droite. Si l'on fait la moyenne de ces différences avec tous les points, on obtient un nombre positif. Dire que la droite est proche du nuage de points veut dire ici que ce nombre est le plus petit possible. Il existe d'autres manières de mesurer la distance entre la droite et les points, et **choisir** une manière de mesurer cette distance est souvent une partie très importante des modèles.

voisins sur lesquels faire la moyenne, mais on pourrait imaginer des modèles plus compliqués avec plus de paramètres).

Autre

On peut imaginer tout un tas d'autres modèles, avec des règles prédéfinies, des paramètres à ajuster, un mélange des deux, ... Ils ont tous des défauts (généralise mal à des nouvelles données, est trop imprécis, ...) et des avantages (est synthétique, apporte un support théorique, ...).

Limites de la modélisation

On observe ici plusieurs points importants, que l'on retrouve dans les Grands Modèles de Langage et les autres modèles génératifs :

- lorsque l'on utilise les données pour calibrer le modèle (trouver la meilleure droite par exemple), la qualité de prédiction du modèle dépend **directement** de la qualité des données (des données éronnées donnent un modèle erroné). La capacité d'un modèle à **généraliser** doit être mesurée, sur des données indépendantes (par exemple les nouvelles données de consommation électrique),
- le **choix** du modèle est important, et influence la prédiction. Certains modèles sont plus adaptés que d'autres dans certains problèmes,
- le modèle nous donne une **approximation** d'un comportement compliqué. La consommation électrique ne dépend probablement pas que de la température! Essayer de prédire la consommation électrique seulement avec la température nous mènera forcément à des approximations, quelque soit le modèle choisit. Les données conditionnent la capacité à modéliser le comportement.

Décomposition **biais-variance-bruit** : en *Machine Learning* (**Apprentissage Automatique**), ces trois derniers points correspondent à la décomposition *biais*, *variance*, *bruit* d'un problème de modélisation basé sur des données. Le *biais* est l'influence des données sur la qualité du modèle, la *variance* correspond au fait de bien modéliser le problème avec ce modèle, et le *bruit* est inhérent à notre problème et à la structure des données. Plus de détails ici : (11).

1.1.2 Apprentissage Automatique (*Machine Learning*), entraînement

Les modèles de **Machine Learning** sont une sous-famille des modèles d'IA. Ce sont des modèles auxquels on va pouvoir apprendre à réaliser des tâches (par exemple, apprendre à déterminer si une photo donnée est une photo de voiture ou non). Souvent, cet apprentissage est effectué en utilisant des **données** (en suivant l'exemple, des milliers de photos dont on sait lesquelles contiennent des voitures). Ces données permettent **d'ajuster** (on parle **d'entraînement**) le modèle pour qu'il **fonctionne bien**.

Dans notre exemple de consommation électrique, la droite est un modèle de Machine Learning! Entraîné avec des données, en choisissant la droite la plus proche du nuage de point. Les paramètres que l'on ajuste ici sont la **hauteur** et la **pente** de la droite. Nous avons calibré le modèle sur une tâche de prédiction de la consommation électrique à partir de la température.

Dans le processus d'entraînement d'un modèle de *Machine Learning*, une partie importante consiste à évaluer le modèle. En effet, **les paramètres** sont choisis pour que le modèle soit le meilleur possible **sur les données d'entraînement**. Par exemple, dans le cas de la droite, on minimise la distance moyenne entre les points et la droite. Si les données sont erronées (erreur de mesure, ...); le modèle va être entraîné à prédire ...des informations erronées (et il sera peut-être très bon pour prédire ces informations). La qualité des données est une partie très importante de ces modèles.

Nombre de paramètres

Là où une droite peut être définie par 2 nombres (sa hauteur et sa pente), les grands modèles de langage récents ont bien plus de paramètres (on compte en dizaines voir centaines de **milliards**). C'est grâce à ces

milliards de paramètres que l'on arrive à modéliser des choses aussi complexes que le langage naturel. C'est de là que vient l'adjectif "Grand" dans « Grand Modèle de Langage ».

L'essor de l'IA ces dernières années (depuis 2010) est dû en partie à des modèles avec un nombre très important de paramètres. Cet essor s'explique car ces modèles permettent de résoudre des tâches complexes (génération de texte, d'image, reconnaissance d'image, jouer et gagner à des jeux complexes comme le go, ...). On parle de **Deep Learning (Apprentissage Profond)** pour désigner une catégorie de modèles (les **réseaux de neurones**) ayant beaucoup de paramètres.

Les idées d'architecture de ces modèles sont plus anciennes que 2010, mais la **quantité de données** amenée par **Internet**, ainsi que le déploiement de **moyens de calcul** plus importants sont les facteurs qui ont permis l'entraînement et la popularisation de ces modèles. Un des points faibles de ces modèles est qu'on ne comprend pas vraiment comment ils fonctionnent. En revanche, on sait observer qu'ils fonctionnent bien ! L'étude de ces modèles est devenue en partie une science expérimentale.

1.2 Quelques exemples de modèles de *Machine Learning*, et les modèles génératifs

Les modèles de Machine Learning fonctionnent pour la plupart en "entrée"- "sortie". Ce sont des programmes qui nous permettent de transformer quelque chose (par exemple une description textuelle d'image), en une autre (une image). Le modèle présenté juste au dessus, utilisé pour prédire la consommation électrique à partir de la température, ne déroge pas à cette règle.

Le modèle YOLO (12) (Figure 4), nous permet de transformer une image, en des "boîtes" annotées : un chien, un vélo, une voiture.

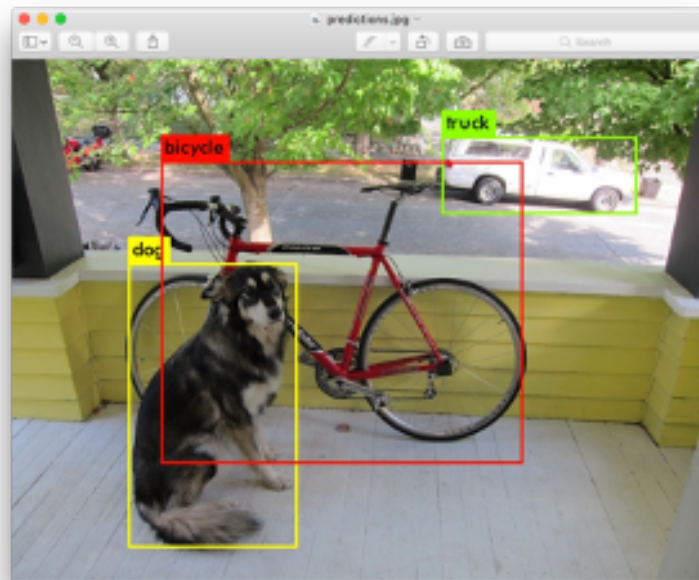


Figure 4 – Un exemple de segmentation d'image avec le modèle YOLO (You Only Look Once) (12)

Le modèle AlphaFold (Figure 5) sert à prédire la structure des protéines en 3 dimensions. En entrée, une séquence d'acides aminés (en pratique c'est simplement une suite de lettres); en sortie, des coordonnées dans l'espace. Ces coordonnées sont celles de la protéine constituée de cette séquence d'acides aminés. La structure en 3 dimensions de la protéine détermine sa fonction (par exemple (13) : transporter de l'oxygène dans le sang).

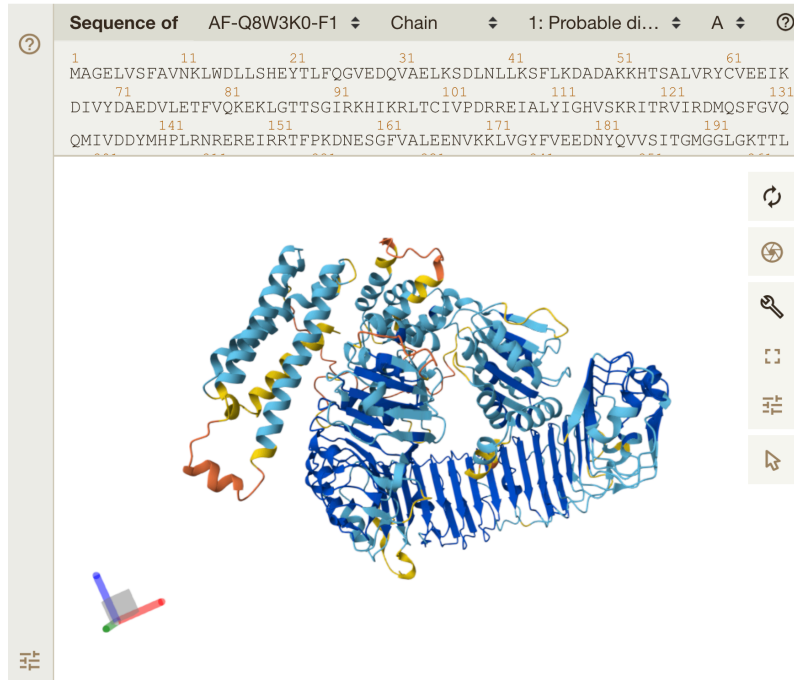


Figure 5 – AlphaFold - Probable disease resistance protein At1g58602 (14)

1.2.1 IA générative

Le terme “IA générative” revient beaucoup dans le débat public. Un modèle d’IA est dit génératif lorsque l’on peut s’en servir pour produire du contenu : **texte, image, vidéo, audio, déplacements dans l’espace**. La plupart du temps, ceux-ci désignent des modèles de Deep Learning. Les LLM, que l’on décrit dans ce document sont des modèles d’IA générative, mais on peut aussi citer Stable Diffusion (à l’origine de l’hippopotame - Figure 1 du début du document).

Attention cependant à l’acronyme **IAG**. Il est fortement ambigu car il désigne plutôt l’IA Générale. L’IA Générale fait référence à un modèle (**fictif**, pour le moment (15)) capable d’effectuer “n’importe quelle tâche cognitive propre aux humains ou autres animaux”². Si on veut absolument utiliser un acronyme pour désigner les IA génératives, on pourra utiliser : IAg.

1.2.2 Et les modèles de langage dans tout ça ?

Les modèles de langage, eux, servent à modéliser ... le langage naturel.

Ce sont des modèles d’Apprentissage Automatiques, entraînés sur des grandes bases de données de texte. En pratique, on peut utiliser ces modèle pour des tâches nécessitant une **compréhension du langage naturel** :

- classifier des phrases (un commentaire donné est-il plutôt *positif* ou *négatif* ?),
- estimer si deux textes sont sémantiquement proches,
- générer du texte,
- ...

Modéliser le langage naturel est loin d’être facile, et on s’en rend compte en constatant le temps qui s’est passé avant l’apparition de robots conversationnels qui fonctionnent bien (environ 2018 avec GPT-2). Il existait avant cela d’autres modèles ; et tout un champ de recherche (*Natural Language Processing*). Ces

2. https://fr.wikipedia.org/wiki/Intelligence_artificielle_g%C3%A9n%C3%A9rale#cite_note-:0-2

modèles étaient cependant bien moins impressionnants que les grands modèles de langage de ces dernières années. Là où prédire la consommation électrique en fonction de la température semble être une tâche accessible, les tâches de traitement du langage présentées plus hauts semblent bien plus complexes. Ne serait-ce qu'encoder la polysémie dans un modèle semble complexe (le même mot a des sens différents selon le contexte).

Mais alors, comment fonctionnent les LLM ?

En fait, modéliser le langage avec des heuristiques et des règles est une tâche trop complexe. Il y a trop de cas particuliers à traiter. L'approche utilisée pour construire un LLM est celle... du Machine Learning :

Au lieu d'expliquer les règles régissant le langage naturel, on construit des modèles qui apprennent à partir d'exemples de phrases.

L'objet de la seconde partie de ce document est donc de présenter plus en détail comment sont construits ses modèles. Quelles sont les données utilisées, et comment on peut leur apprendre des tâches.

Néanmoins, juste avant d'ouvrir le capot et de présenter les modèles de langages, nous nous attardons sur les impacts de ces derniers. En effet, du fait de leur taille et de leurs capacités, ces modèles présentent de **nombreux risques** ; qu'il convient d'expliquer. L'usage de ces modèles ne se fait pas sans coûts, loin de là.

1.3 Les impacts des grands modèles génératifs

Le développement et l'utilisation des grands modèles que l'on vient de décrire a des **impacts désastreux**, sur plusieurs plans. Des employés exploités pour visualiser des contenus toxiques (16) aux mini-centrales nucléaires produites à la chaîne pour alimenter des centres de données (17), le paysage n'est pas très beau à voir...

Pour une introduction à ces impacts, nous conseillons l'article et la présentation vidéo des JRES (Journées Réseaux de l'Enseignement Supérieur), qui présentent très bien le sujet : (18). Citons également l'association PauselA : (15) qui produit de nombreux contenus sur le sujet. Pour des études plus approfondies, on pourra regarder l'institut de recherche Epoch AI (19).

Pour rentrer dans les détails, les modèles d'apprentissage profonds posent de sérieuses questions, et ce tout au long de leur cycle de vie. De la fabrication des composants électronique à l'utilisation finale des modèles, en passant par leur entraînement et leur mise à disposition ; chacune de ces étapes apporte son lot de limitations et de problématiques.

1.3.1 Fabrication des composants et durée de vie du matériel

La production des composants électroniques nécessaires à l'entraînement et à la mise à disposition de ces modèles nécessite des matières premières rares (on parle de terres rares). Ces terres rares sont parfois extraites dans des conditions désastreuses. Le tantale par exemple, utilisé dans les condensateurs ; est principalement produit en République Démocratique du Congo (43% en 2017 (20)). L'article de radio-canada (21) décrit une mine où l'extraction se fait manuellement, parfois par des enfants.

Une fois produites, les cartes graphiques ont *a priori* une durée de vie très courte (1 à 3 ans selon (22)), due à un taux d'utilisation très élevé. L'article cite un architecte de Google, selon qui il serait possible de prolonger cette durée - au doux prix de les utiliser moins souvent et donc d'amortir moins vite leur achat.

La question du recyclage des composants n'est pas du tout résolue ; alors que les cartes graphiques contiennent des produits dangereux pour l'environnement (23). Le rapport sur la durabilité de la compagnie Nvidia (24) (ayant un quasi-monopole sur les cartes graphiques utilisées pour l'entraînement des grands modèles de langage) reste bizarrement très timide sur le sujet ; malgré leur sublime phrase d'introduction : "Green isn't just our corporate color" (23).

1.3.2 Nettoyage et annotation des données

L'entraînement des modèles d'apprentissage automatique nécessite beaucoup de données. La qualité des données impactant directement la qualité des modèles, les entreprises qui développent ces derniers ont parfois recours à de méthodes de nettoyage et d'étiquetage des données **manuelles**. C'est le cas par exemple de l'entreprise OpenAI (développant ChatGPT) qui, en 2022, a sous-traité des tâches d'annotations d'images à une entreprise au Kenya (Sama). Les employés étaient payés autour de 2\$ de l'heure ; pour visionner des contenus toxiques : abus sexuel d'enfants, suicide, torture, ... Une enquête sur le sujet a été effectuée par le TIME magazine (16).

Des questions intéressantes sont à régler autour des sources de données et des droits d'auteurs. Il semble à peu près de notoriété publique que les modèles de langage aient été entraînés avec des données soumises au droit d'auteur (25). Enfin, l'existence même du droit d'auteur sur les contenus générés avec ces modèles pose également question : (26).

1.3.3 Entraînement des modèles, mise à disposition

Les ordinateurs utilisés dans les phases d'entraînement et d'exploitation des modèles ont besoin d'électricité pour fonctionner, et de systèmes de refroidissements pour ne pas surchauffer. Ces consommations électriques sont loin d'être négligeables, et en sérieuse hausse (27).

Même si l'accent est mis sur les émissions CO₂³, il ne faut pas oublier que ce n'est pas la seule manière de mesurer l'impact environnemental. On peut également calculer la dépletion des ressources abiotiques (diminution de la quantité de matières premières non-vivantes), ou la consommation d'énergie primaire pour produire les composants et les acheminer à bon port. De manière générale, les impacts environnementaux liés à la fabrication sont rarements calculés.

Les modèles génératifs sont tellement gros, que même une fois entraînés, le fait de transformer une entrée en sortie nécessite beaucoup de puissance de calcul, et donc d'électricité, d'eau, ... La consommation en électricité des centres de calculs est tellement importante que Google envisage même d'y dédier des centrales nucléaires (17). Pour avoir un ordre de grandeur des consommations électriques des modèles (hors entraînement), on pourra consulter l'article (29). Citons également cette étude (27), présentant des prédictions de consommation électriques dans les 5 prochaines années, pour l'entraînement de ces modèles aux Etats-Unis. En bref, la croissance exponentielle en besoin de calculs pour l'entraînement des modèles risque d'être stoppée par une production électrique insuffisante, d'ici 2030.

En plus de l'électricité, les centres de calculs ont besoin d'eau pour refroidir les composants (30). Ceci pose des questions liées aux conflits d'usages d'une ressource vitale.

1.3.4 Impacts sociétaux

Les impacts sociétaux liés à ces modèles sont également importants. On peut citer (15) pour une liste de ces dangers ; qui incluent : le chômage de masse, la perte de compétence (et donc une irréversibilité de l'utilisation de ces modèles), désinformation, deepfakes, ...

Pour conclure cette partie, voici quelques (contre)-arguments sur l'adoption de ces technologies au quotidien : (31).

1.4 *Facultatif - Zoom sur les réseaux de neurones*

Les réseaux de neurones sont une famille de modèles de **Machine Learning**. Ils sont utilisés dans quasiment tous les grands modèles récents. Grossièrement, les réseaux de neurones sont découpés en couches. Pour transformer une entrée en sortie, les réseaux de neurones vont transformer successivement l'entrée

3. Les papiers de recherches décrivant les phases d'entraînement des nouveaux modèles de langage diffusent parfois des estimations de CO₂. Meta estime par exemple dans un papier que le développement d'une génération de modèles à émis autour de 1000 tonnes équivalent CO₂ (28).

avec chacune des couches. Chaque couche a ses paramètres (à la manière de notre droite), qui peuvent être réglés pour minimiser l'erreur du modèle sur des données d'entraînement. On parle de réseaux profonds (*Deep Learning*) lorsqu'ils ont beaucoup de couches.

Le succès des réseaux de neurones dans les modèles récents s'explique par plusieurs points :

- l'algorithme d'optimisation des paramètres de ces réseaux est très souple (méthode du gradient) ; et très bien implémenté (Automatic Differentiation (32))
- l'architecture en couche s'adapte à beaucoup de cas, et bénéficie d'une communauté *open-source* active (avec par exemple PyTorch (33)), ce qui favorise la réutilisation et le développement de nouveaux modèles,
- les réseaux de neurones fonctionnent bien pour modéliser beaucoup de phénomènes qu'on arrive pas à expliciter avec des règles (générer des images, du texte, ...),
- les recherches sur ces modèles ont été (et sont toujours) accélérées par les financements privés (globalement toutes les grosses entreprises du web : les GAFAM⁴ et autres entreprises américaines, des entreprises chinoises comme Alibaba, ...), qui ont leurs propres laboratoires de recherche.
- ...

2 Explication du fonctionnement des LLM

Un LLM (*Large Language Model*) est un modèle de *Machine Learning* qui peut être utilisé pour des tâches nécessitant une compréhension du langage naturel. Par exemple, mais pas que, pour générer du texte. Comme on l'a décrit précédemment, ces modèles ont pour la plupart plusieurs dizaines de milliards de paramètres. La phase d'entraînement d'un tel modèle est très coûteuse en infrastructure⁵ - elle demande beaucoup de matériel informatique - et en temps - les phases d'entraînement peuvent durer plusieurs mois.

Les robots conversationnels (*chatbots*) comme **ChatGPT** sont construits au-dessus de modèles appelés **modèles de fondations**. Afin de bien comprendre comment ces *chatbots* fonctionnent, on commence par présenter les **modèles de fondation**. On peut les voir un peu comme des **moteurs de voiture**. En suivant cette comparaison, les agents conversationnels peuvent être vus comme les différentes voitures que l'on peut construire autour d'un moteur. Pour se repérer, l'agent conversationnel ChatGPT est construit au-dessus de modèles de fondations appelés GPT (*Generative Pre-trained Transformer*). Plusieurs modèles existent : GPT-2, GPT-3, GPT-4...

Le terme *modèle de fondation* n'est pas propre aux modèles de langage ; et désigne en général un *grand* modèle dont l'entraînement est conséquent ; mais qui peut être **spécialisé** pour de nombreuses tâches différentes : génération de texte, calcul de similarité de deux phrases, agent conversationnel, ...

2.1 Modèle de fondation, *tokens* et machines à camemberts

Les modèles de fondation sont le **noyau** des agents conversationnels comme *ChatGPT*. Le fonctionnement précis de ces modèles varie. Aussi, on présente ici les grandes idées des modèles déjà existants, mais il faut bien garder à l'esprit que ces domaines évoluent très vite.

On a vu dans la première partie de ce document que la plupart des modèles de *Machine Learning* permettent de transformer une "entrée" en une "sortie". Les modèles de langage ne font pas exception. Une grande part de ces modèles sont entraînés sur la tâche suivante :

- en entrée, une séquence de mots,
- en sortie, prédire la suite **probable** de la phrase.

4. Désignait les entreprises américaines : Google, Amazon, Facebook, Apple, Microsoft. Aujourd'hui, certains noms ont changé : Google est devenu Alphabet et Facebook Meta. L'entreprise OpenAI (qui développe ChatGPT) a été rachetée par Microsoft.

5. Le modèle V3 de l'entreprise chinoise DeepSeek a *a priori* (c'est encore sujet à débat, à prendre avec précaution) nécessité bien moins de puissance de calcul que ses homologues pour leur entraînement, pour des tailles et des performances similaires (34). Ceci va peut-être rebattre les cartes sur les moyens de calculs nécessaires à l'entraînement de ces modèles.

D'une certaine manière, à l'usage, ces modèles se comportent comme le mécanisme d'autocomplétion des téléphones. Lorsqu'un LLM est utilisé pour générer des lignes et des lignes de textes, c'est tout simplement que la sortie est réinjectée dans l'entrée. On parle de mécanisme auto-régressif (*auto-regressive* ou *causal* en anglais). D'autres modèles de langage sont entraînés sur une tâche de remplissage de textes à trous (c'est le cas de BERT (35)). L'entrée des modèles est très souvent appelée "**prompt**" (en français, invite ou **instruction**).

Plus concrètement, étant donnée une séquence de mots (l'entrée); la première étape de l'algorithme de génération de texte consiste à la convertir en *tokens* (**jetons** en Français). Les *tokens* sont l'unité de langage la plus fine utilisée dans les *LLMs* (6). Ils sont de taille comparable aux mots (environ 3-4 lettres), mais il faut absolument garder à l'esprit que le modèle est utilisé pour générer des *tokens* et non des mots voire des lettres. La sortie du modèle est un *token* probable suivant la séquence donnée en entrée. Pour résumer :

- "entrée" = "séquence de *tokens*",
- "sortie" = "un *token* probable suivant la séquence d'entrée".

Le découpage d'un texte en *token* est effectué à l'aide d'un *tokenizer*. C'est un algorithme assez simple (comparé au modèle de génération), qui transforme une séquence de lettres en une séquence de *tokens*. Notons que ce découpage en *token* ne contient absolument aucune **sémantique**. C'est uniquement un découpage statistique des textes. Les mêmes *tokens* sont utilisés pour générer des textes dans différentes langues. Du point de vue du LLM, les langages n'ont pas de différences : tout n'est que texte, et séquence de *tokens*.




Figure 6 – Découpage d'une phrase en *tokens* (36)

On peut faire la liste de tous les *tokens* possibles. C'est ce qu'on appelle le "vocabulaire" du modèle. La taille du vocabulaire varie selon les modèles; en gros de 30 000 à plus de 100 000. Le découpage en *tokens* à l'avantage par rapport à un découpage en mots de ne pas créer des vocabulaires trop grands (imaginez la taille d'un vocabulaire constitué de tous les mots de la langue française, puis de l'allemand, anglais, ...). Et il a l'avantage par rapport à un découpage en caractères de ne pas créer des séquences encodées trop longues. C'est une sorte de *juste milieu* entre ces 2 découpages (4).

Rentrons encore un peu plus dans le détail. À l'étape précédente la sortie, il est généré un **camembert** - ou *diagramme circulaire*. Ce camembert est en fait la sortie la plus utile du modèle⁶. Il y est représenté les différents *tokens* pouvant apparaître à la suite de la séquence fournie en entrée. Seulement, la subtilité des LLM réside dans le fait qu'à chacun des *tokens* est assigné une probabilité d'apparition. Certains *tokens* ont plus de chance d'apparaître à la suite de cette séquence (Figure 8), et sont donc représentés avec une plus grande part de camembert. Le modèle a été entraîné à prédire avec le plus de précision possible ce camembert - et est donc en quelque sorte une machine à camembert.

Le point **central** de la génération de texte par les LLM est le suivant : le prochain *token* est **tiré au hasard**. Ce sont des modèles de nature **probabiliste**. Deux **entrées** identiques peuvent produire en pratique 2 sorties différentes. Le prochain *token* est tiré au hasard, oui, mais **pas uniformément**. Il est tiré au hasard, selon les probabilités d'apparition du camembert : imaginez comme une roue de la fortune, qui tourne et s'arrête sur une part du camembert. Les parts les plus grosses (= les *tokens* les plus probables) auront plus de chances d'être tirés au sort.

Plus précisément, différentes méthodes existent pour tirer au sort le prochain *token*; chacune se basant sur le camembert. La méthode peut être choisie lors de la génération de texte, par les utilisatrices et utilisateurs des modèles. Citons en vrac différentes méthodes, directement issues de l'article (38) :

- *Greedy* (avare) : peut-être la moins utile, elle consiste à choisir le *token* avec la plus forte probabilité; tuant le côté aléatoire du modèle;

6. C'est la sortie la plus utile car elle contient plus d'information que simplement le *token* suivant. Ce camembert va pouvoir être utilisé pour spécialiser le modèle de fondation dans tout un tas de tâches variées.

token_id	token
0	!
1	"
2	#
3	\$
4	%
...	...
128251	< reserved_special_token_243 >
128252	< reserved_special_token_244 >
128253	< reserved_special_token_245 >
128254	< reserved_special_token_246 >
128255	< reserved_special_token_247 >

128256 rows × 1 columns

Figure 7 – Vocabulaire du modèle Llama 3.2 1B

Suite de la phrase : 'La souris de l'ordinateur est un' - Llama 3.2 1B

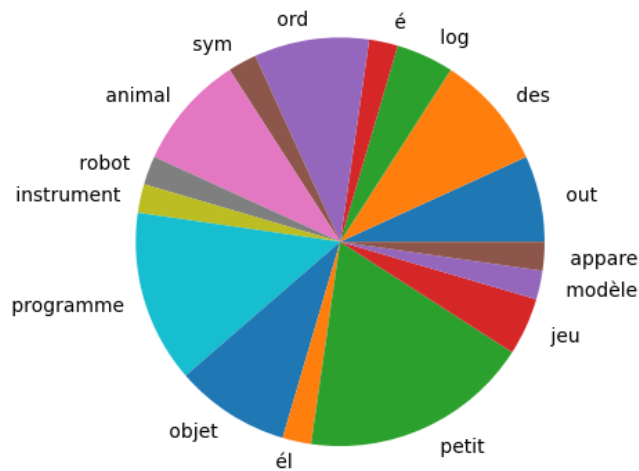


Figure 8 – Prédiction du prochain token - Llama 3.2 1B (37)

- *Beam search* (recherche par faisceau) : un peu plus évoluée, cette méthode consiste à regarder une étape plus loin. Elle choisit le couple des 2 *tokens* suivants ayant le plus de chance d’apparaître. Par exemple, pour la séquence d’entrée : “La souris”; on observe 2 suites possibles : “La souris de l” et “La souris est un”. Dans le premier cas, le *token* “de” avait une probabilité de 0.8, puis “l” une probabilité de 0.3. Dans le second, le *token* “est” avait une probabilité de 0.6 et “un” de 0.9. Au total, la première séquence a donc une probabilité $0.8 \times 0.3 = 0.24$ d’apparaître, et la seconde $0.6 \times 0.9 = 0.54$. L’algorithme implémentant le modèle nous fournit comme *token* de sortie la seconde option (la plus probable).
- *Sampling* (échantillonnage) : ici, chaque *token* peut être choisi, en fonction de sa probabilité : les *tokens* avec des plus fortes probabilité ont plus de chance d’être tirés. On peut imaginer qu’on fait tourner la roue (celle donnant les probabilités du prochain *token*); et qu’on choisit le *token* là où la roue s’arrête. C’est le mode classique décrit juste au-dessus.
- *Top K* : même méthode que l’échantillonnage, mais on ne garde que les K *tokens* les plus probables.
- *Top p* : de manière un peu similaire à la précédente, cette méthode consiste à ne garder que les *tokens* les plus probables, pour que la somme de leurs probabilités dépasse p . On choisit ensuite le prochain *token* au hasard, comme dans la méthode “échantillonnage”.

L’étape consistant à utiliser le modèle pour prédire un *token* est appelée **l’inférence**.

Dans l’exemple ci-dessus, le camembert présente le *token* le plus probable comme étant le *token* “petit”. Néanmoins, cela ne garantit pas que la suite de la phrase sera “petit”, un autre *token* pourrait être tiré au hasard. Par exemple, le mot “programme”, “objet” ou “animal” peuvent aussi apparaître. On observe ici la polysémie du mot “souris”; tous les cas de figure sont envisagés. C’est en partie permis par le côté “probabiliste” du modèle.

Pour prédire des textes plus longs, il suffit simplement de redonner au modèle la séquence que nous lui avons fait compléter (d’où le terme **auto-régressif**). Faisons cela sur l’exemple ci-dessus; la sortie obtenue est la phrase : “La souris de l’ordinateur est un objet qui permet de manipuler les informations sur un ordinateur.”. Point important : sur les dernières prédictions de la phrase (Figure 9) les probabilités sont nettement plus franches en faveur d’une souris d’ordinateur. Ceci donne un exemple de l’importance du **contexte** : en donnant plus de précisions aux agents conversationnels, vous augmentez vos chances d’obtenir du texte généré cohérent par rapport à votre question.

Hallucination

Suite de la phrase :
'La souris de l'ordinateur est un objet qui permet de manipuler les informations sur un'
Llama 3.2 1B

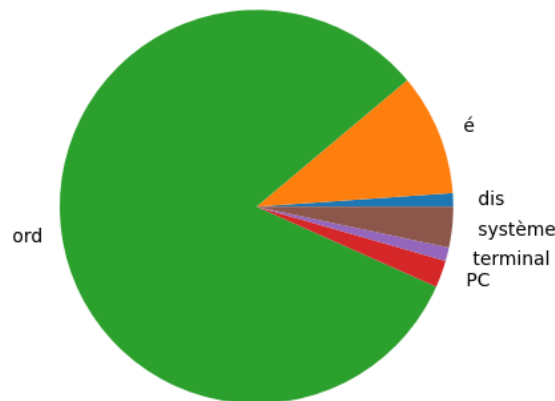


Figure 9 – Prédiction du prochain token - Llama 3.2 1B (37)

Le LLM est donc simplement entraîné sur une tâche de prédiction d'un *token* suivant une séquence donnée. Il n'y a là dedans **absolument** aucune notion de ce qui est **vrai** ou **faux**. Il n'y a que des probabilités, qui penchent en faveur de générations de contenus cohérents. Ceci laisse la porte **grande ouverte** au fait que ces modèles produisent des contenus **faux, offensants, dangereux**, ... On désigne ce comportement par le terme "hallucination".

Les hallucinations ne sont pas étonnantes à ce stade de l'entraînement, la notion de vérité peut être très loin de celle de *token* probable. Cependant, elle persiste jusqu'aux agents conversationnels ; et donc il ne faut absolument pas perdre de vue que toute information sortant d'un tel modèle **doit être vérifiée et sourcée**. Ainsi, toute information que l'on n'est pas capable de vérifier - ou de faire vérifier - doit être vue comme **inutilisable**. Pour contourner ces problèmes, les **RAG** (Retrieval Augmented Generation) sont une excellente parade : ils consistent à utiliser les LLM pour chercher dans une base de données ; puis à fournir la source de leur réponse. En bref, même si aucune notion de *vérité* n'est câblée dans les modèles, et que les contenus générés avec de tels outils sont potentiellement faux ; cela ne signifie pas qu'on ne peut pas les utiliser comme interface de langage naturel. En effet, ce sont tout de même des modèles que l'on peut utiliser pour analyser des requêtes d'utilisateurs, par exemple.

Température

La température est un nombre que les utilisateurs et utilisatrices des LLM peuvent choisir pour modifier la génération de texte. Ce nombre représente à quel point les *tokens* générés sont "inattendus". En théorie, n'importe quel nombre positif peut être accepté, mais en pratique, les outils permettant d'utiliser les modèles restreignent aux nombres entre 0 et 1 (certains acceptent jusqu'à 2). La température est donc un coefficient qui vient modifier les probabilités d'apparition du prochain *token*, juste avant le camembert. Cependant, l'ordre des *tokens*, de celui ayant la plus faible probabilité d'apparition à celui avec la plus élevée, ne change pas avec la température⁷. Voir Figure 10 pour quelques exemples.

Une température proche de 0 va pénaliser les *tokens* avec des faibles probabilités d'apparition et favoriser ceux avec des fortes probabilités d'apparition : cela augmente le manichéisme du texte généré. Une température de 0 rend le modèle complètement déterministe (identique au comportement *greedy* de génération

7. Pour que l'ordre des *tokens* change, il faudrait mettre une température négative. Dans ce cas, les *tokens* avec les plus faibles probabilités d'apparitions deviennent ceux avec les plus fortes, et *vice-versa*. Pas très utile... ça explique pourquoi les outils restreignent aux températures positives.

de texte), le *token* avec la plus haute probabilité d'apparition est **systématiquement** généré. Au contraire, une température plus éloignée de 0 va favoriser les *tokens* avec des faibles probabilités d'apparition, et pénaliser ceux avec de fortes probabilités d'apparition. En fait, avec une température au-delà de 2, des probabilités tellement faibles sont assignées aux *tokens* les plus probables que le texte généré n'a plus aucun sens (c'est une sorte de créativité extrême qui consiste à raconter littéralement n'importe quoi).

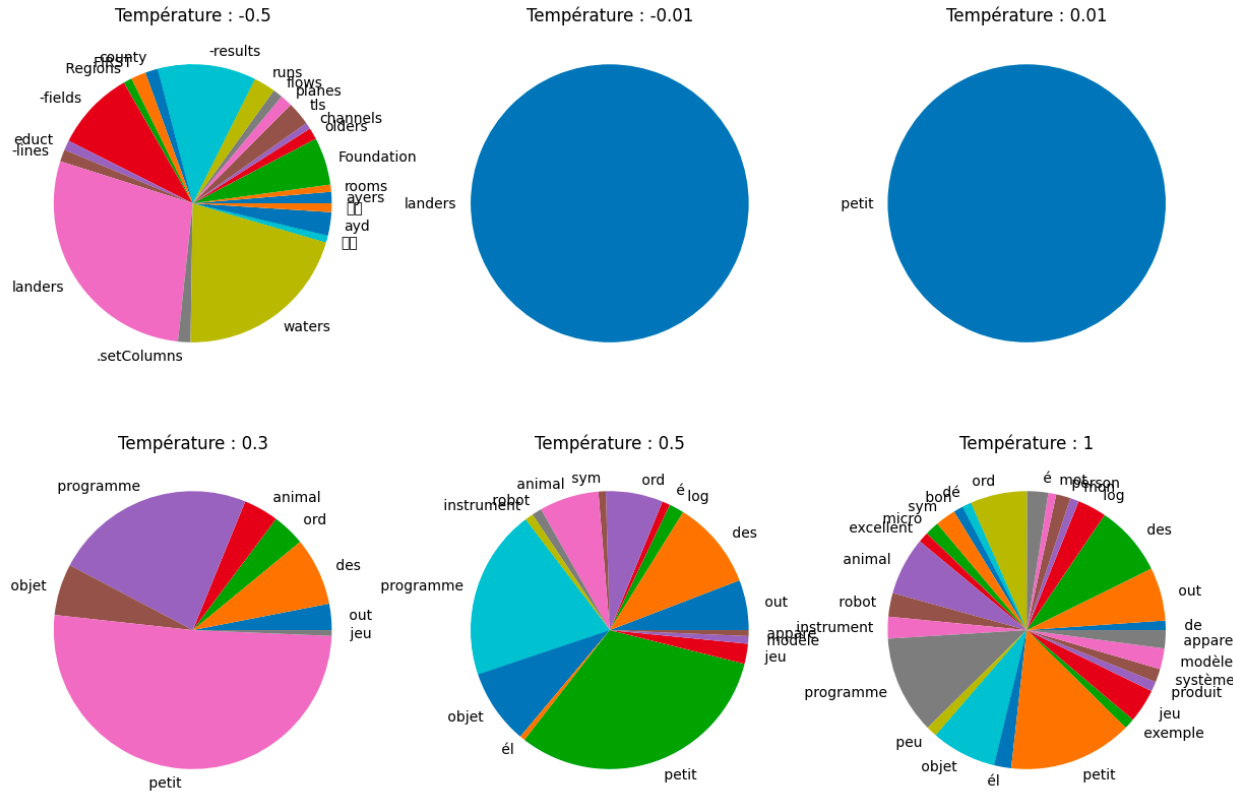


Figure 10 – Suite de la phrase : “La souris de l’ordinateur est un”, avec différentes températures - Llama 3.2 1B

2.2 Entraînement des modèles de fondation

Les grands modèles de langage sont des modèles d'apprentissage automatique. On leur apprend la structure du langage à partir de données textuelles. De la même manière que la droite de l'exemple initial est composée de deux paramètres, ces modèles de fondations sont composés de milliards - dizaines / centaines de milliards - de paramètres.

Les paramètres ne sont pas agencés n'importe comment, loin de là. Le fait que ces modèles fonctionnent pour les tâches de traitement du langage naturel est dû à leur architecture. Comme présenté rapidement dans la partie sur les réseaux de neurones; ces modèles sont des modèles de **Deep Learning**; constitués d'une succession de différentes couches de neurones. La petite révolution sur l'architecture des **LLM** est due à des couches dites “d'attention”, basée sur des “*transformers*”. Initialement présentée dans un article (39) de chercheurs de Google, cette technologie est utilisée dans quasiment tous les LLM aujourd'hui.

Si on initialise le modèle que l'on a présenté juste au dessus avec des paramètres quelconques (mais en gardant l'architecture du modèle bien sûr!), il prédit tout simplement n'importe quoi! Dans la Figure 11 on présente le camembert produit par le modèle avec des paramètres quelconques. Pour information, en essayant

de continuer la génération, nous avons obtenu la séquence suivante : “La souris de l’ordinateur estmpjeses-per Annie ?Compar ?? clothes_sink QByteArray eitherenciónats Sustainabilityimanpyn ?? compatibility nth either”. Clairement, le modèle n’est pas capable de produire du texte intelligible.

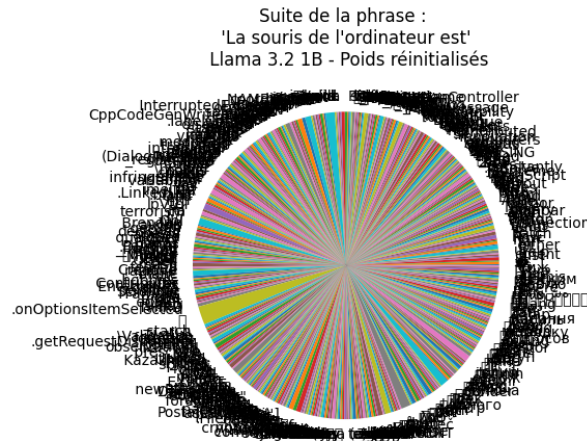


Figure 11 – Prédiction du prochain token avec des paramètres réinitialisés - Llama 3.2 1B (37)

L’entraînement du modèle de fondation consiste alors à trouver les meilleures valeurs possibles pour les paramètres, afin que les camemberts générés produisent du langage cohérent (dans la partie précédente, les phrases sont produites avec le modèle entraîné bien sûr).

Comme dans le modèle de prédiction de température, la phase d’entraînement s’appuie sur des **données**. Et ça tombe bien, Internet contient énormément de contenus textuels ! En pratique, de grands jeux de données sont constitués à partir de données venant d’Internet : Wikipédia (40), GitHub (base de données de projets de programmation) (41), ArXiv (base de donnée d’articles scientifiques) (42), ... On pourra consulter les Tables 1 et 2 de l’article (8) pour une vue d’ensemble des différents jeux de données utilisés pour l’entraînement des modèles. On peut également citer le jeu de données FineWeb (43) ; qui reconstruit des données semblables à celles de l’entraînement de certains *LLM*. Les données brutes sortant d’Internet ne sont pas utilisées directement dans l’entraînement. Un énorme travail de nettoyage de données (filtrage de contenus indésirables, gestions des contenus de différentes langues, déduplication, suppression d’informations personnelles, ...) est fait en amont.

Pour les modèles Llama, les données nettoyées constituent moins de 5TB de stockage disque (!) (28). C’est très peu, mais ces données ne sont pas organisées, et il est potentiellement impossible de retrouver facilement une information dedans. Ces modèles sont alors entraînés avec environ 1 billion de *tokens* (mille milliards, 1 000 000 000 000). Une fois entraîné, on peut voir ces modèles comme une synthèse des données d’entraînement. L’utilisation de ces modèles pour générer du texte ne va pas directement chercher dans les bases de données d’entraînement. Il est plus correct de dire que les données d’entraînement sont synthétisées dans le modèle.

L’entraînement se déroule de la manière suivante :

- On prend une partie du texte d’entraînement, découpé en *tokens*. Par exemple, 9 tokens consécutifs dans la page Wikipédia du gratin dauphinois : “Le”, “grat”, “in”, “dau”, “ph”, “inois”, “est”, “un”, “plat”.
- On tronque cette séquence en ne gardant que les 8 premiers : “Le”, “grat”, “in”, “dau”, “ph”, “inois”, “est”, “un”.
- On utilise le modèle pour générer une distribution de probabilité pour le prochain *token* (c’est-à-dire, pour chaque *token* du vocabulaire, on estime sa probabilité d’apparition à la suite de la séquence des 8 *tokens*). Cela nous donne un camembert, comme précédemment.
- Le camembert obtenu est comparé avec le résultat attendu : le *token* “plat”. Très concrètement, une

distance est calculée entre le camembert avec un seul *token* (" plat"), et celui généré.

- Les paramètres du modèle sont modifiés en fonction de cette comparaison. La distance dépend de tous les paramètres du modèle. L'algorithme d'entraînement consiste alors à mettre à jour les paramètres pour qu'au prochain coup, la distance soit plus petite.

Cette dernière étape est de loin la plus coûteuse en calcul. Elle est faite avec un algorithme (rétro-propagation de gradient) qui permet de calculer à quel point chaque paramètre doit être modifié pour améliorer la sortie du modèle. C'est loin d'être facile, d'autant plus que le nombre de paramètres est tellement grand que de nombreux problèmes se posent sur le chemin. Les implémentations de cet algorithme sont très astucieuses (auto-differentiation (32), (33)).

En pratique, l'algorithme d'entraînement du modèle génère en une passe le camembert de chaque portion de la séquence d'entrée :

- le camembert pour l'entrée : "Le",
- le camembert pour l'entrée : "Le", " grat",
- ...
- le camembert pour l'entrée : "Le", " grat", "in", " dau", "ph", "inois", " est", " un".

Ceci nous fait 8 camemberts (pour cet exemple) à comparer. Les paramètres sont alors mis à jour en une seule fois pour les 8 camemberts : la distance utilisée pour mettre à jour le modèle est la moyenne des distances sur les 8 camemberts.

2.3 Du modèle de fondation à l'agent conversationnel

Le modèle de fondation, comme son nom l'indique, est à la racine de beaucoup d'applications différentes. C'est un modèle que l'on utilise pour des tâches nécessitant de comprendre le langage naturel ou d'en générer. C'est donc une très bonne base pour produire tout un tas de modèles plus spécifiques.

On peut par exemple spécialiser le modèle de fondation pour obtenir un algorithme de classification de sentiment : savoir si une phrase est positive ou négative⁸ ; ou si un texte contient des données personnelles ou non.

Une autre application serait de construire un programme qui nous calcule la similarité entre deux phrases⁹. Ici, on présente une de ces spécialisations, la création d'un **agent conversationnel**. C'est celle qui a vraiment popularisé les LLM, notamment avec ChatGPT. Cependant, il ne faut pas oublier que, selon les usages, l'agent conversationnel n'est pas toujours la meilleure solution.

Afin d'obtenir ces agents conversationnels, plusieurs méthodes ont été imaginées par les développeurs et développeuses des modèles. Les différentes méthodes utilisées dépendent des organismes produisant les modèles, et ne sont pas toujours expliquées (il y a parfois une sorte de "secret de fabrication"). Ainsi, il convient de rester **prudent** sur toute cette partie : plus particulièrement que pour le reste du document, cette section ne contient pas de "vérité générale" qui s'applique et s'appliquera à tous les modèles ; mais plutôt quelques grandes idées que l'on retrouve dans plusieurs modèles (ceux qui ont publié des informations sur leur processus d'entraînement). On cite la *survey* récent (44) qui présente de manière très claire les différentes méthodes pouvant être utilisées pour spécialiser un modèle de fondation.

Voici un exemple de complétion d'une séquence de *tokens* avec un modèle de fondation (Llama 3.2 1B) :

Entrée :

"J'ai dans mon placard des tomates pelées, du riz et du poulet. Qu'est-ce que je peux cuisiner ?"

Sortie du modèle de fondation (Llama 3.2 1B) :

"Je suis contente de pouvoir cuisiner! J'ai un petit gars qui a besoin d'attention et de nourriture. Je vais le nourrir."

8. <https://huggingface.co/tabularisai/multilingual-sentiment-analysis>

9. https://www.sbert.net/docs/sentence_transformer/training_overview.html

Et voici la sortie de l'agent conversationnel :

```
Sortie de l'agent conversationnel (Llama 3.2 1B Instruct) :  
"C'est une question classique! Voici quelques idées de plats que vous pouvez  
préparer avec ces ingrédients :  
__Riz et poulet aux tomates__ :  
Faites chauffer du riz et du poulet dans une poêle, ajoutez des tomates pelées  
coupées en dés et du sel. Servez chaud."
```

Ces deux réponses sont complètement différentes. Dans la première, le texte généré est intelligible et cohérent avec l'entrée, mais ne répond pas du tout à l'attente de l'utilisateur (en plus d'expliquer un biais sexiste). Le second en revanche, y répond très bien. Ce problème de "bien répondre aux attentes de l'utilisateur" est fréquemment appelé : **problème d'alignement**. C'est la différence entre ce qu'on attend d'un modèle d'apprentissage automatique, et ce qu'il sort vraiment de ce modèle. Le modèle de fondation est entraîné à générer un *token* probable suivant une séquence de *token*, pas du tout à répondre aux questions d'un utilisateur.

Nous allons expliquer comment on peut arriver à un tel agent.

Dans la littérature, le terme **pre-training** (pré-entraînement) est utilisé pour parler de l'entraînement du modèle de fondation, et **post-training** (post-entraînement) pour parler de la spécialisation du modèle de fondation en agent conversationnel¹⁰. Le pré-entraînement, que l'on a décrit juste avant, est coûteux (plusieurs mois, plusieurs millions - au grand minimum - de dollars d'investissement matériel). Le post-entraînement nécessite *a priori* bien moins de temps. Une fois le modèle complètement entraîné (ou même seulement avec le modèle pré-entraîné), il est possible de continuer à le spécialiser avec les méthodes présentées ici, pour qu'il intègre des données propres à un domaine. Cela permettra d'obtenir un modèle pour son propre cas d'usage (par exemple, retenir des acronymes, ou devenir spécialiste de la peinture flamande du 19e siècle). On parle dans ce cas de *fine-tuning* (ajustement)¹¹.

On trouve deux types de méthodes majoritaires dans le post-entraînement.

- Ajustement supervisé (*Supervised Fine Tuning (SFT)*),
- Ajustement par renforcement (*Reinforcement Learning (RL)*).

2.3.1 Ajustement supervisé

Cette phase consiste à continuer l'entraînement que l'on a décrit précédemment (du modèle de fondation), mais avec un nouveau jeu de données. Au lieu de prendre des suites de *tokens* issues d'Internet, on utilise des jeux de données de conversations entre un agent (aimable, prêt à aider les gens, ...) et un utilisateur. À l'origine, ces jeux de données ont été constitués par des **experts humains**, payés pour répondre à des questions, afin que les échanges soient enregistrés dans le jeu de données. Désormais, il est très probable que ces jeux de données soient partiellement produits par des LLM (les experts humains peuvent être aidés dans la rédaction de leur réponse par des LLM, par exemple).

Ci-dessous, on montre un exemple de conversation, dont les réponses ont été remplies par des humains. On observe qu'une même question admet plusieurs formulations de réponse différentes.

C'est aussi l'occasion de conditionner le modèle à formuler des réponses spécifiques dans des cas précis : dire qu'il ne connaît pas la réponse à une question, éviter de répondre aux questions sensibles, ... L'objectif est de transformer le modèle de fondation, qui est une sorte de synthèse de tout internet, en une synthèse des différents experts ayant participé à rédiger les conversations d'entraînement.

La procédure d'entraînement est la même que pour le modèle de fondation (prédiction du camembert pour le prochain *token*, et ajustement des paramètres en fonction du *token* réel); mais il y a une nouveauté

10. Avec les termes *pre-training* et *post-training*, on se demande à quel moment se passe le "training". En fait, l'entraînement est la combinaison du *pre-training* et du *post-training*.

11. <https://github.com/meta-llama/llama-cookbook>

text string	role string	lang string
Combien de livres Hunger Games existe-t-il ? ça fait environ combien de pages ?	prompter	fr
Il existe 4 livres Hunger Games: -Hunger Games (384 pages) -Hunger Games: L'embrasement (400 pages) -...	assistant	fr
Il y a 3 ouvrages principaux à la série Hunger Games, avec un préquel. En version française, ça...	assistant	fr
Il y a 3 livres principaux (1 prequel) et ça fait beaucoup de pages : Livre 1 : 399 Livre 2 :400 Livr...	assistant	fr

Figure 12 – Question - Réponse entre un utilisateur et un assistant. OASST1 (45) - Hugging Face

importante : le rajout de certains *tokens* spéciaux. Ci-dessous, on a affiché l’encodage en *tokens* d’un échange utilisateur-assistant, juste avant qu’il soit donné au modèle pour l’entraînement.

```
<|im_start|>user<|im_sep|>Combien de livres Hunger Games existe-t-il ? ça fait environ combien de pages ?<|im_end|><|im_start|>assistant<|im_sep|>Il existe 4 livres Hunger Games: -Hunger Games (384 pages) -Hunger Games: L'embrasement (400 pages) -Hunger Games: La Révolte (339 pages) -Hunger Games: La Ballade du Serpent de l'Oiseau Chanteur (1712 pages) 384+400+339+1712= 2835 pages (Le nombre de page est calculé à partir des éditions originales de ces livres, et non des versions françaises)<|im_end|><|im_start|>assistant<|im_sep|>
```

Figure 13 – Encodage d’un échange utilisateur-assistant (45) (46)

Pour entraîner le modèle à générer du texte similaire à une conversation entre un assistant et un utilisateur, des *tokens* spéciaux sont rajoutés au vocabulaire du modèle. Par exemple, écrivons ces tokens : “<|im-start|>”, “assistant”, “<|im-stop|>”. Ces *tokens* sont différents des textes : “<|im-start|>”, “assistant” et “<|im-stop|>” tokenisés ! Ce sont vraiment des nouveaux *tokens*, avec des identifiants de *tokens* différents dans le vocabulaire (voir Figure 7). Ils sont ensuite intégrés artificiellement aux données d’entraînement et, leur sens est **appris**. Les textes générés par le modèle ressembleront statistiquement plus souvent aux réponses d’un assistant aimable après les occurrences de ces *tokens* spéciaux de conversation.

Il existe d’autres *tokens* spéciaux ; par exemple pour utiliser des outils (4). En rajoutant un token encodant le fait d’aller chercher une information sur internet, il est possible de conditionner le modèle à produire ce token lorsque la réponse générée est floue. Le serveur donnant accès au modèle peut alors détecter ce token et insérer dans le *prompt* le résultat d’une recherche internet. Un autre exemple est un *token* signalant le début de l’exécution de code, généré par le modèle dès qu’il semble pertinent d’utiliser du code pour répondre aux questions utilisateurs.

Dans tous ces cas, la sémantique de ces *tokens* est **apprise**, à travers tous les exemples des données d’entraînement.

2.3.2 Ajustement par renforcement

L’objectif de l’ajustement par renforcement d’entraîner le modèle en lui faisant générer plusieurs réponses et de **juger** lesquelles sont les meilleures. Ce type de méthode permet de **renforcer** l’apparition de certains

type de réponses à l'avenir.

Cependant, lorsqu'une question n'admet pas de bonne ou de mauvaise réponse, par exemple : "Fais moi une blague sur les ordinateurs", il est complexe de juger de la qualité d'une réponse. Une solution utilisée (notamment pour ChatGPT 3 (47)) est d'entraîner un **autre** modèle à prédire un score de qualité sur les réponses du LLM. Ces modèles sont appelés des *reward models* (modèles de récompenses); et sont entraînés avec des **retours humains**¹², pour donner des scores ressemblant à ce qu'un humain prédirait. En utilisant ce modèle auxiliaire, on peut dire au modèle de langage quelles réponses sont les plus appropriées parmi celles qu'il nous a données. Il a été observé que cela développe sa capacité à se comporter comme un humain (47).

3 Comment représenter des *tokens* avec des nombres ? Les *embeddings*

Dans cette partie, on présente comment les *tokens* sont utilisés dans les modèles. C'est l'étape d'*embedding*, qui désigne la transformation des *tokens* en des séquences de nombres (**vecteurs**).

3.1 Vecteurs

On commence par expliquer brièvement ce qu'est un **vecteur** dans notre contexte. C'est en pratique tout simplement une liste de nombres. Par exemple :

$$\begin{bmatrix} 0 \\ 1.2 \\ 0.2 \end{bmatrix}$$

est un vecteur. La longueur de la liste (ici 3) est appelée **dimension du vecteur**. Chacun des éléments du vecteur est appelé une "composante", ou "coordonnée". Un point crucial est qu'il est possible d'effectuer des opérations sur des vecteurs. On définit ainsi la somme de 2 vecteurs comme un nouveau vecteur, dont les composantes sont la somme des composantes des 2 vecteurs. Par exemple :

$$\begin{bmatrix} 0 \\ 1.2 \\ 0.2 \end{bmatrix} + \begin{bmatrix} 18 \\ 2.1 \\ 0.9 \end{bmatrix} = \begin{bmatrix} 18 \\ 3.3 \\ 1.1 \end{bmatrix}.$$

Un vecteur à 2 dimensions peut être **représenté** dans un plan, et à 3 dimensions dans **l'espace**. Il est commode de les représenter avec une flèche. Prenons l'exemple (Figure 14) des deux vecteurs :

$$\begin{bmatrix} 3 \\ -2 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \end{bmatrix}.$$

Une raison pour visualiser les vecteurs par une flèche issue de l'origine, qui pointe sur les coordonnées du vecteur, est que ça permet de bien comprendre ce que représente la somme de 2 vecteurs. Le vecteur obtenu en sommant les deux premiers correspond simplement à la composition des deux flèches. Il faut imaginer qu'un vecteur encode un déplacement, et que la somme encode la composition des déplacements.

Si le vecteur a plus de 3 dimensions, on peut garder l'intuition qu'il correspond à une flèche, mais dans un espace que l'on ne peut pas trop visualiser ! Malheureusement, c'est très souvent le cas...

12. Le terme précis : RLHF, *Reinforcement Learning with Human Feedback*.

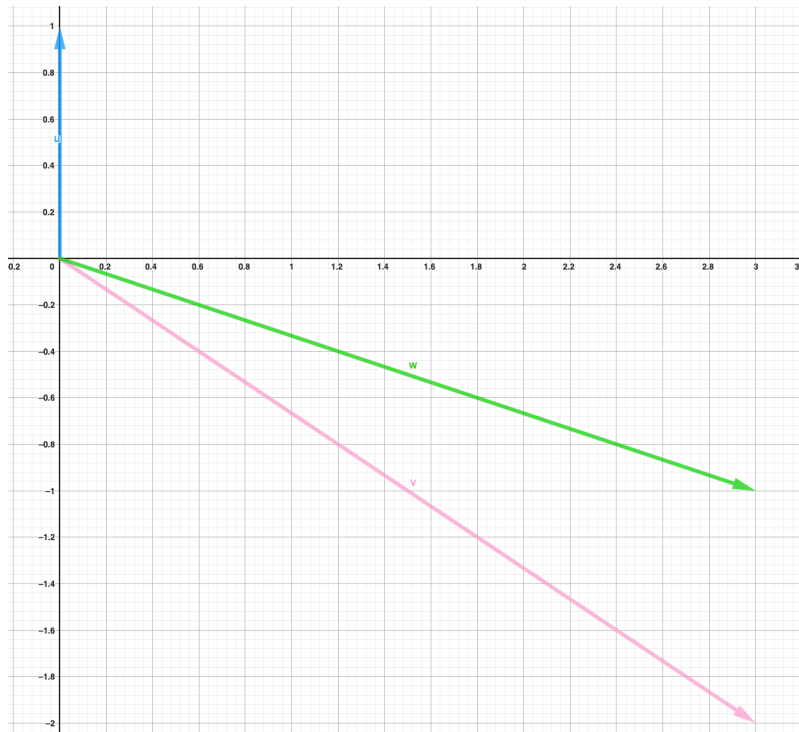


Figure 14 – Le vecteur vert est la somme du rose et du bleu. Il pointe sur l’endroit où on atterrit après avoir suivi le rose, puis le bleu.

3.2 Embedding

En traitement du langage, l’*embedding* correspond à la transformation du texte en **vecteurs**. La dimension des vecteurs varie, mais peut être de plusieurs centaines - jusqu’à plusieurs milliers ! Difficile alors de se représenter ces vecteurs. Cette représentation des *tokens* est utile car la traduction *tokens-vecteurs* est faite afin que 2 séquences de tokens sémantiquement proches (par exemple “chocolat” et “cacao”) aient des vecteurs dont les directions sont proches.

Les usages sont très nombreux et ne s’arrêtent pas à la génération de texte : calcul de similarité de 2 phrases, classification selon le sens d’un texte, recherche documentaire, ... Selon l’usage, les vecteurs dont on a besoin peuvent être radicalement différents. On pourra consulter le site SBERT¹³, présentant des implémentations.

Attention, le terme *embedding* désigne à la fois l’action de transformer un *token* en un vecteur, et le vecteur obtenu.

On recense plusieurs méthodes pour transformer des *tokens* (ou directement des mots) en vecteurs. Certaines existaient avant la popularisation des LLM.

Par exemple Word2Vec (48) ; issu de chez Google en 2013. Word2Vec désigne deux modèles de *Machine Learning* ; que l’on peut utiliser pour transformer des mots en vecteurs. Pour jouer avec ces représentations, vous pouvez regarder **Cemantix**¹⁴, un jeu qui consiste tout simplement à deviner un mot, et qui s’appuie justement sur les distances entre vecteurs, obtenues avec Word2Vec. Il est intéressant de remarquer que ces deux modèles sont entraînés sur une autre tâche, et que la représentation des mots en vecteur est en quelque sorte une étape intermédiaire qui permet aux modèles de bien réaliser cette tâche. Le premier modèle est entraîné à remplir un texte avec un unique mot manquant, au centre (un peu comme les LLM

13. Sentence Transform (SBERT) : <https://sbert.net/index.html>.

14. Cemantix : <https://cemantix.certitudes.org/>.

finalement) ; et le second à trouver les mots en entourant un autre (48). Comme il est complexe de déterminer si un *embedding* est correct ; on est un peu obligé de passer par ce genre de chemin détourné pour entraîner des modèles d'*embeddings*. Remarquons qu'une fois les *embeddings* calculés pour tous les mots d'un vocabulaire, la traduction mot-vecteur est très rapide (il suffit de regarder dans le dictionnaire!). Dans le même style, on peut citer GloVe (49) ; de la même époque.

Dans le cas des LLM, les *tokens* sont transformés en vecteur dès leur entrée dans le modèle. Les composantes de chaque vecteur d'*embedding* du vocabulaire sont initialisées au hasard, puis **appries** lors de l'entraînement : ce sont des paramètres du modèle. En un certain sens, comme avec Word2Vec, les *embeddings* sont une conséquence de l'entraînement du LLM à prédire le prochain *token* d'une séquence donnée. Pour mieux prendre en compte les informations sur la position de chaque *token* dans une phrase, les LLM ont en plus des "*positional embeddings*" (vecteurs de position). Ce sont d'autres vecteurs, appris lors de l'entraînement ou non, qui encodent la position d'un *token* dans une phrase. Là encore, plusieurs variantes existent (*positional embedding* (50), *rotary positional embedding* (51) (52) (53)).

Une des petites révolution des LLM consiste à utiliser une architecture appelée "transformer" (50), qui permet aux *embeddings* de chacun des *tokens* (sous forme de vecteur) d'être modifiés pour prendre en compte les autres *tokens* de la séquence. En utilisant des *transformers*, il est donc possible de faire évoluer les *embeddings* initiaux (appris pendant l'entraînement) pour qu'ils deviennent plus aboutis (deux *embeddings* du même mot seront alors différents selon le contexte). Ceci donne lieu aux modèles d'*embeddings* modernes ; dans lesquels l'étape de traduction *token*-vecteur doit être refaite à chaque fois que l'on veut transformer un *token* en vecteur dans un contexte précis. Impossible dans ce cas de créer un dictionnaire *tokens*-vecteurs ! Ces modèles sont alors plus coûteux que les anciens (Word2Vec par exemple) à utiliser, car les *transformer* restent des architectures avec beaucoup de paramètres.

Récemment, un retour en arrière est observé : les *embeddings statiques* (54). Ils fonctionnent lors de la traduction comme les modèles style Word2Vec : ce sont des dictionnaires *tokens*-vecteurs. Un exemple est Model2Vec (55). C'est une technique permettant d'utiliser des modèles prenant en compte le contexte (*transformer*) pour créer un dictionnaire. Une fois créé, le dictionnaire ne permet pas de prendre en compte le contexte, et l'*embedding* d'un *token* est toujours le même. Néanmoins, la technique est suffisamment au point pour que cela n'affecte pas trop la qualité des *embeddings*.

4 Écosystème et historique des LLM

Les grands modèles de langage ont profondément transformé le domaine du traitement automatique du langage grâce à l'augmentation importante des ressources informatiques, à la disponibilité accrue de grandes bases de données textuelles et à l'innovation majeure des architectures Transformer. Cette partie explique clairement ce qui différencie les LLM et retrace leur évolution technologique et commerciale.

4.1 Différencier les LLM

Les modèles de langage se distinguent notamment selon leur accès, leur taille, leur contexte et leur application.

Ouverture des modèles

Bien que de nombreux modèles de langage se revendiquent ouverts ou « open source », leur transparence réelle varie considérablement. L'ouverture d'un modèle implique plusieurs critères : disponibilité du code source, accès aux données et poids d'entraînement, existence de licences explicites, ainsi que qualité et exhaustivité de la documentation associée notamment concernant la phase d'entraînement du modèle.

Il apparaît que certains modèles offrent une transparence significative sur la plupart de ces critères (par exemple, OLMo (56) d'Allen AI, qui publie ses données d'entraînement, son code source et des documentation détaillées), alors que d'autres modèles populaires demeurent essentiellement opaques (DeepSeek

R1 (57), dont les données d'entraînement restent fermées, ou Llama 3.3 (58), qui nécessite une licence spécifique et dont les données ne sont pas divulguées).

Ainsi, le degré d'ouverture influence directement la possibilité d'utiliser ces modèles dans des contextes scientifiques, éducatifs ou responsables.

Des ressources comme Opening Up ChatGPT (59) et l'Open Source AI Index (60) proposent des évaluations complètes et à jour du degré d'ouverture des modèles de langage. On pourra également consulter le tableau Figure 15.

Accès aux modèles

Les **modèles fermés** comme GPT-4 (61) ou Claude (62) sont accessibles uniquement *via* des interfaces propriétaires (API¹⁵ ou applications web), sans possibilité d'accéder à leur code ou à leurs paramètres internes. Cela facilite leur utilisation mais limite la transparence. Les **modèles à tendance open-source** comme LLaMA (63) offrent un accès complet à leur code source et à leurs paramètres, permettant ainsi une utilisation plus flexible, notamment pour des développements spécifiques.

Taille et nombre de paramètres

La taille d'un LLM impacte directement ses capacités et ses besoins en ressources. Les modèles **petits** (quelques millions à quelques centaines de millions de paramètres) fonctionnent facilement sur des ordinateurs standards et conviennent pour des tâches simples.

Les modèles **grands** (plusieurs dizaines à centaines de milliards de paramètres) nécessitent une infrastructure importante, et sont adaptés pour des tâches complexes comme la traduction ou la génération détaillée de texte.

Taille du contexte

Le contexte est la quantité de texte (en nombre de *tokens*) que le modèle peut traiter en une seule fois. Des modèles comme Gemini 1.5 (64) traitent des contextes allant jusqu'à un million de *tokens*, permettant ainsi de gérer des textes ou documents très longs.

Domaine d'application

Certains modèles sont spécialisés dans des domaines précis (médical, juridique, financier), offrant une grande précision mais limitée à ce secteur. D'autres modèles restent généralistes, permettant une utilisation polyvalente.

4.2 Historique des LLM

L'évolution des LLM peut être divisée en plusieurs grandes périodes historiques.

4.2.1 1990-2010 : Approches statistiques et neuronales

Les premiers modèles de langage s'appuyaient sur des n-grammes (65) (séquences de n mots consécutifs, permettant d'estimer la probabilité d'apparition des mots) (66). Ils exploitaient les statistiques de co-occurrence de mots, mais étaient limités aux séquences courtes en raison d'une complexité exponentielle (67). Vers la fin des années 1990, l'utilisation de réseaux neuronaux récurrents (RNN (68) - *Recurrent Neural Network*) – c'est-à-dire des modèles qui traitent l'information séquentiellement en conservant une mémoire des éléments précédents – a été introduite (68). Ces approches ont rapidement évolué avec l'émergence

15. API : *Application Programming Interface*, désigne ici une manière automatisée (machine à machine) d'accéder au modèle. Ceci permet par exemple de créer un programme informatique qui utilise de manière automatique des modèles hébergés ailleurs que sur son ordinateur.

Project	Availability						Documentation					Access		
	Open code	LLM data	LLM weights	RL data	RL weights	License	Code	Architecture	Preprint	Paper	Modelcard	Datasheet	Package	API
OLMo 7B Instruct	✓	✓	✓	✓	✓	✓	✓	✓	✓	X	✓	✓	✓	~
BLOOMZ	✓	✓	✓	✓	~	~	✓	✓	✓	✓	✓	✓	X	✓
AmberChat	✓	✓	✓	✓	✓	✓	~	~	✓	X	~	~	X	✓
Open Assistant	✓	✓	✓	✓	X	✓	✓	~	~	X	X	X	✓	✓
OpenChat 3.5 7B	✓	X	✓	X	✓	✓	~	✓	✓	~	~	X	✓	~
Pythia-Chat-Base-7...	✓	✓	✓	✓	X	✓	✓	~	~	X	~	~	✓	X
Cerebras GPT 111...	~	✓	✓	✓	✓	~	X	✓	~	X	X	✓	X	✓
RedPajama-INCITE...	~	✓	✓	✓	✓	~	~	~	X	X	✓	✓	X	~
dolly	✓	✓	✓	✓	X	✓	✓	✓	~	X	X	X	✓	X
Tulu V2 DPO 70B	✓	X	~	✓	✓	~	~	~	~	X	~	~	~	✓
MPT-30B Instruct	✓	~	✓	~	X	✓	✓	~	X	X	~	X	✓	~
MPT-7B Instruct	✓	~	✓	~	X	✓	✓	~	X	X	✓	X	✓	X
trix	✓	✓	✓	~	X	✓	✓	~	X	X	X	X	~	✓
Vicuna 13B v 1.3	✓	~	✓	X	X	~	✓	X	✓	X	~	X	✓	~
minChatGPT	✓	✓	✓	~	X	✓	✓	~	X	X	X	X	✓	✓
ChatRWKV	✓	~	✓	X	X	✓	~	~	~	X	X	X	✓	~
BELLE	✓	~	~	~	~	X	~	✓	✓	X	X	~	X	X
WizardLM 13B v1.2	~	X	~	✓	✓	~	~	✓	✓	X	X	X	X	X
Airoboros L2 70B G...	~	X	~	✓	✓	~	~	~	X	X	~	~	X	X
ChatGLM-6B	~	~	✓	X	X	✓	~	~	X	~	X	X	X	✓
Mistral 7B-Instruct	~	X	✓	X	~	✓	X	~	~	X	X	X	~	✓
WizardLM-7B	~	~	X	✓	~	~	~	✓	✓	X	X	X	X	X
Qwen 1.5	~	X	✓	X	✓	X	~	~	X	X	X	X	~	✓
StableVicuna-13B	~	X	~	~	~	~	~	~	~	X	~	X	X	~
Falcon-40B-instruct	X	~	✓	~	X	✓	X	~	~	X	~	X	X	X
UltraLM	X	X	~	✓	~	X	X	~	✓	X	~	~	X	X
Yi 34B Chat	~	X	✓	X	✓	~	X	X	✓	X	X	X	X	~
Koala 13B	✓	~	~	~	X	~	~	~	X	X	X	X	X	X
Mixtral 8x7B Instruct	X	X	✓	X	~	✓	X	~	~	X	X	X	~	X
Stable Beluga 2	X	X	~	X	✓	~	X	~	~	X	~	X	X	~
Stanford Alpaca	✓	X	~	~	~	X	~	✓	~	X	X	X	X	X
Falcon-180B-chat	X	~	~	~	~	X	X	~	~	X	~	X	X	X
Orca 2	X	X	~	X	✓	X	X	~	~	X	~	X	X	~
Command R+	X	X	X	✓	✓	~	X	X	X	X	~	X	X	X
Gemma 7B Instruct	~	X	~	X	~	X	X	~	X	X	✓	X	X	X
LLaMA2 Chat	X	X	~	X	~	X	X	~	~	X	~	X	X	~
Nanbeige2-Chat	✓	X	X	X	✓	~	X	X	X	X	X	X	X	~
Llama 3 Instruct	X	X	~	X	~	X	X	~	X	X	~	X	X	~
Solar 70B	X	X	~	X	~	X	X	X	X	X	~	X	X	~
Xwin-LM	X	X	~	X	X	X	X	X	X	X	X	X	X	~
ChatGPT	X	X	X	X	X	X	X	X	~	X	X	X	X	X

Figure 2: Openness of 40 text generators described as open, with OpenAI's ChatGPT (bottom) as closed reference point. Every cell records a three-level openness judgement (✓ open, ~ partial or X closed). The table is sorted by cumulative openness, where ✓ is 1, ~ is 0.5 and X is 0 points. RL may refer to RLHF or other forms of fine-tuning aimed at fostering instruction-following behaviour. For the latest updates see: <https://opening-up-chatgpt.github.io>

Figure 15 – Tableau recensant le degré d'ouverture des grands modèles de langage entraînés par instructions - (59)

des LSTM (69) (Long Short-Term Memory, des réseaux capables de conserver l'information sur de longues périodes grâce à des mécanismes de mémoire) et, plus tard, des GRU (70) (Gated Recurrent Units) une variante simplifiée des LSTM utilisant des portes pour contrôler le flux d'information (69).

Les réseaux neuronaux récurrents comme les LSTM présentaient des difficultés d'entraînement sur de grands corpus, notamment à cause de la complexité des calculs nécessaires. L'arrivée des réseaux neuronaux à propagation avant (*feed-forward*), comme le Neural Network Language Model (NNLM) (71), a introduit les représentations distribuées des mots (*word embeddings*). Ces représentations vectorielles ont permis de saisir les similarités sémantiques et les relations complexes entre les mots, dépassant les limitations des approches statistiques traditionnelles basées sur les n-grammes.

Cette période est également liée à l'amélioration des techniques de lissage (pour attribuer des probabilités non nulles aux séquences jamais observées dans les données d'entraînement) pour les modèles n-grammes (65), et leur optimisation pour des corpus de plusieurs milliards de mots (72). Ensuite, Mikolov et al. (73) ont commencé à explorer les réseaux de neurones récurrents pour la modélisation du langage, introduisant notamment des architectures RNN optimisées qui réduisaient considérablement le temps d'entraînement tout en améliorant les performances prédictives. Ces travaux ont posé les fondations conceptuelles et techniques des futures architectures comme les GRU (74), en démontrant la viabilité des approches neuronales récurrentes pour capturer des dépendances à long terme dans le langage.

4.2.2 2011-2016 : Word Embeddings et avancées RNN

Word2Vec (48) et GloVe (49) ont introduit des méthodes pour convertir les mots en vecteurs denses, dans lesquels la proximité entre les vecteurs reflète la similarité sémantique des mots. Ces techniques, en permettant de placer les mots dans un espace numérique à dimension réduite, ont amélioré l'efficacité des réseaux neuronaux récurrents. La diffusion de bibliothèques d'apprentissage profond, tel que TensorFlow, a aussi accéléré ces progrès. Durant cette période, le mécanisme d'attention, procédé qui permet de pondérer l'importance de chaque élément d'une séquence lors de la prédiction en calculant des scores de pertinence, a commencé à être exploité dans des applications telles que la traduction automatique (74).

Avant ces innovations, les mots étaient souvent représentés sous forme "*one-hot*" (un vecteur avec un seul 1 à l'index du mot et des 0 ailleurs), ce qui ne permettait pas de capturer la similitude entre les mots et engendrait des vecteurs de très grande dimension. Les méthodes de *word embeddings* comme Word2Vec ou GloVe viennent plutôt fournir une représentation compacte et informative.

Avancées des RNN et mécanisme d'attention

L'intégration des word embeddings a amélioré les représentations d'entrée des réseaux de neurones récurrents (RNN), ce qui a permis d'optimiser leurs performances pour des tâches comme la classification de textes ou la traduction automatique. L'arrivée de frameworks tels que Theano, TensorFlow et PyTorch (33) a facilité l'expérimentation et la mise en œuvre de modèles plus avancés, comme les architectures LSTM (Long Short-Term Memory) et GRU (Gated Recurrent Units).

Cependant, ces modèles éprouvaient des difficultés à gérer des dépendances lointaines. Ils peinaient à se souvenir des informations antérieures et à relier des mots éloignés dans le texte. L'introduction du mécanisme d'attention, associée à la montée en puissance du calcul parallèle sur GPU, a ouvert la voie aux architectures sans récurrence, annonçant l'émergence des Transformers.

Le mécanisme d'attention (39) permet au modèle d'attribuer un score de pertinence à chaque mot d'une séquence. Contrairement aux méthodes antérieures, qui réduisaient l'ensemble de la séquence à un vecteur fixe, cette approche crée un vecteur contextuel différent pour chaque étape de la prédiction, en fonction de la position du mot. Ce procédé améliore la prise en compte des relations entre mots éloignés et, par conséquent, la qualité des prédictions.

4.2.3 2017-2019 : Apparition des Transformers

L'introduction de l'architecture Transformer par Vaswani et al. (50), présentée dans « Attention is All You Need », représente une rupture par rapport aux approches séquentielles des RNN. Le Transformer exploite le mécanisme d'attention multi-têtes pour analyser simultanément toutes les positions d'une séquence, réduisant ainsi le temps d'entraînement et capturant les dépendances à longue distance. Le self-attention calcule pour chaque mot un score de pertinence par rapport aux autres, permettant de générer un vecteur contextuel adapté à chaque étape de prédiction. Cette structure encodeur-décodeur a ouvert la voie à des modèles de plus en plus performants, comme GPT-1 et BERT en 2018, suivis par GPT-2, XLNet et RoBERTa en 2019, qui ont tous contribué à améliorer les performances sur diverses tâches.

Modèles phares de cette période

GPT-1 (Generative Pre-trained Transformer 1) (75) : GPT-1 a été l'un des premiers modèles à démontrer l'efficacité du pré-entraînement massif d'un Transformer (en utilisant uniquement la partie décodeur) pour la génération de texte et le transfert vers d'autres tâches. Il utilise 117 millions de paramètres.

BERT (Bidirectional Encoder Representations from Transformers) (35) : BERT utilise la partie encodeur du Transformer et est pré-entraîné sur deux tâches : la prédiction de mots masqués ("Masked Language Modeling") et la prédiction de la phrase suivante ("Next Sentence Prediction"). Cette approche bidirectionnelle (considérant à la fois le contexte gauche et le contexte droit) permet à BERT de capturer des représentations contextuelles très riches. BERT existe en différentes tailles, la plus grande ayant 340 millions de paramètres (BERT-Large).

GPT-2 (76) : GPT-2 est une version beaucoup plus grande de GPT-1 (jusqu'à 1,5 milliard de paramètres). Il a été entraîné sur un corpus encore plus vaste et a démontré des capacités de génération de texte impressionnantes, parfois difficiles à distinguer d'un texte écrit par un humain.

XLNet (77) : XLNet est un modèle qui combine les avantages de l'approche autorégressive (comme GPT) et de l'approche bidirectionnelle (comme BERT). Il utilise une méthode de pré-entraînement appelée "Permutation Language Modeling" qui lui permet de prendre en compte toutes les permutations possibles de l'ordre des mots dans une phrase.

RoBERTa (Robustly Optimized BERT Pretraining Approach) (78) : RoBERTa est une version optimisée de BERT. Les auteurs ont montré qu'en modifiant certains hyperparamètres, en entraînant le modèle plus longtemps et sur plus de données, et en supprimant la tâche de "Next Sentence Prediction", ils pouvaient améliorer significativement les performances de BERT.

4.2.4 2020-2021 : Montée en puissance et adoption

Entre 2020 et 2021, les modèles de langage de très grande échelle émergent, qui, en combinant augmentation de taille et apprentissage en contexte, ont transformé l'approche des tâches de traitement du langage.

GPT-3 (79) (développé par OpenAI) a repoussé les limites de la taille des modèles en intégrant 175 milliards de paramètres, soit plus de 100 fois la taille des modèles précédents comme BERT-Large, et plus de dix fois la taille de GPT-2. L'apparition de ce modèle introduit des capacités inédites telles que le *few-shot learning* (voir ci-dessous), sans nécessiter un entraînement spécifique pour chaque tâche. Pour la première fois, un modèle pouvait s'adapter dynamiquement à de nouvelles tâches simplement à travers des instructions textuelles et quelques exemples dans le prompt, là où les approches précédentes nécessitaient systématiquement un *fine-tuning* avec des milliers d'exemples étiquetés pour chaque nouvelle application.

- **Zero-Shot** : Le modèle est capable d'effectuer une tâche sans avoir été explicitement entraîné sur des exemples de cette tâche. On lui donne simplement une description de la tâche en langage naturel (par exemple, "Traduire de l'anglais au français : cheese =>") et il génère la sortie attendue ("fromage").
- **One-Shot** : Le modèle reçoit un seul exemple de la tâche à accomplir (par exemple, "Traduire de l'anglais au français : sea otter => loutre de mer, cheese =>") avant de générer la sortie ("fromage").

- **Few-Shot** : Le modèle reçoit quelques exemples (typiquement entre 2 et une dizaine, voire une centaine) de la tâche (par exemple, plusieurs paires de mots anglais-français) avant d'effectuer la tâche sur une nouvelle entrée. L'idée est de montrer au modèle le format attendu de la sortie, ainsi que quelques exemples de la relation entrée-sortie.
- **In-context learning** : L'ensemble de ces méthodes (zero-shot, one-shot et few-shot) sont des formes de "in-context learning". Le modèle apprend à effectuer une tâche à partir d'instructions ou d'exemples fournis dans le contexte de la requête, sans modification de ses paramètres (contrairement au fine-tuning).

Ces approches réduisent considérablement, voire éliminent, le besoin de grands ensembles de données étiquetées pour chaque nouvelle tâche. Cela rend les modèles de langage beaucoup plus polyvalents et faciles à adapter à de nouveaux usages tout en se rapprochant de la capacité humaine à apprendre rapidement à partir de quelques exemples. Des modèles comme T5 (Text-to-Text Transfer Transformer) (80) et Megatron-LM (81) ont poursuivi la tendance à l'augmentation de la taille et des capacités.

Diversification et spécialisation des modèles

- T5 (Text-to-Text Transfer Transformer) (80) développé par Google, est un modèle unifié qui reformule toutes les tâches de NLP comme des tâches de texte à texte. Par exemple, la traduction est formulée comme "Translate English to German : That is good. =>", la classification de sentiment comme "Classify the sentiment of this sentence as positive or negative : I love this movie. =>", etc. Cette approche permet de pré-entraîner un seul modèle sur un mélange de tâches et de le fine-tuner facilement sur de nouvelles tâches. T5 a été pré-entraîné sur le corpus "Colossal Clean Crawled Corpus" (C4).
- Megatron-LM (81) (et développements ultérieurs) : Megatron-LM est un projet de NVIDIA visant à entraîner de très grands modèles de langage Transformer. Il se concentre sur les techniques permettant de distribuer l'entraînement sur de nombreux GPU, car les modèles de cette taille ne peuvent pas tenir sur un seul appareil.
- Codex (OpenAI) : Codex est un exemple de spécialisation des modèles de langage. Il est basé sur GPT-3, mais il est fine-tuné sur un grand corpus de code source provenant de dépôts publics (comme GitHub). Codex est capable de générer du code à partir de descriptions en langage naturel, d'aider à compléter du code, de traduire entre différents langages de programmation, et d'expliquer le fonctionnement du code. Il alimentait des outils comme GitHub Copilot.

4.2.5 2022 : Démocratisation et open-source

Fin 2022, ChatGPT, basé sur une version améliorée de GPT-3 (souvent appelée GPT-3.5 ou gpt-3.5-turbo), a permis de rendre les grands modèles de langage accessibles au grand public via une interface web intuitive. Contrairement aux versions antérieures, disponibles uniquement via une API payante pour les développeurs, ChatGPT offre une interaction en langage naturel qui ne requiert aucune compétence technique particulière. Les réponses, à la fois fluides et pertinentes – capables d'expliquer des concepts complexes, de rédiger des textes dans divers styles ou même de générer du code – ont suscité un véritable effet viral sur les réseaux sociaux, transformant ainsi l'usage des LLM en une expérience quotidienne.

2022 a également marqué l'accélération du mouvement open-source des LLM. Des projets comme BLOOM, développé par le collectif BigScience (rassemblant plus de 1000 chercheurs et entraîné sur le supercalculateur Jean Zay), offrent un modèle multilingue dont le code, les données d'entraînement et les poids sont entièrement accessibles, permettant de générer du texte dans 46 langues et 13 langages de programmation. De leur côté, les modèles tels qu'OPT de Meta AI – bien que les versions les plus volumineuses soient soumises à des licences restrictives – ainsi que LLaMA et StableLM (82), illustrent comment l'open-source contribue à démocratiser l'accès, accélère la recherche, et renforce la transparence des modèles en facilitant l'audit des systèmes.

4.2.6 2023-2024 : Multimodalité et spécialisation

À partir de 2023, les LLM intègrent la multimodalité, c'est-à-dire la capacité de traiter et générer simultanément du texte et des images. Par exemple, GPT-4, notamment dans sa version multimodale souvent

désignée GPT-4V (pour « Vision »), peut analyser des images, répondre à des questions sur leur contenu ou générer des légendes explicatives. D'autres modèles, tels que Flamingo (83) de DeepMind ou PaLM-E (84) de Google, combinent des capacités linguistiques avec des encodeurs visuels pour étendre leur application à des domaines comme la robotique, tandis que LLaVA (85) émerge comme initiative open-source.

La concurrence s'est également intensifiée avec l'arrivée de nouveaux acteurs et l'émergence de modèles spécialisés. Des systèmes comme Claude d'Anthropic (62), Gemini de Google (64) et LLaMA de Meta (28), proposés en différentes tailles et souvent accompagnés de versions optimisées pour le dialogue, offrent des fenêtres de contexte considérablement étendues – jusqu'à 1 million de tokens pour Gemini 1.5 Pro (64) ou 200k pour Claude 3 (62) – permettant de traiter des textes plus longs et complexes. Des modèles dédiés à des tâches spécifiques, tels que ceux pour la traduction automatique, la génération de code (comme Codex (86) ou Code LLaMA (87)), la rédaction scientifique, l'analyse financière ou encore la réponse à des questions médicales (par exemple, Med-PaLM (88) de Google (88) et BioBERT (89)), ainsi que dans le domaine juridique avec LegalBERT (90), indiquent une spécialisation croissante des LLMs.

4.3 Accès et déploiement des LLM

Différentes méthodes existent pour utiliser ces modèles selon les besoins :

Interfaces en ligne :

Mistral, ChatGPT, Claude.ai ou Gemini proposent des interfaces faciles d'accès via navigateur ou API.

- Mistral : <https://chat.mistral.ai/chat>,
- ChatGPT : <https://chatgpt.com/>,
- Claude.ai : <https://claude.ai>,
- Gemini : <https://gemini.google.com/?hl=fr>.

Déploiement local :

Des outils simples (Ollama (91), LM Studio (92), llama.cpp (93)) permettent d'utiliser les modèles localement pour garantir confidentialité et contrôle.

- Ollama : <https://ollama.com/search>,
- LM Studio : <https://lmstudio.ai/models>,
- llama.cpp : <https://github.com/ggml-org/llama.cpp>.

API pour développeurs :

OpenAI, Anthropic, Google Gemini et Mistral AI proposent des API pour intégrer facilement ces modèles dans des applications.

- API de OpenAI : <https://platform.openai.com/docs/overview>,
- API de Claude par Anthropic : <https://www.anthropic.com/api>,
- API de Google Gemini : <https://ai.google.dev/>,
- API de Mistral AI : <https://docs.mistral.ai/api/>.

Plateformes de modèles :

HuggingFace et Replicate (94) facilitent la comparaison et l'évaluation des performances des différents modèles disponibles.

- HuggingFace spaces : <https://huggingface.co/spaces>,
- Replicate : <https://replicate.com/explore>.

L'écosystème des LLM est en constante évolution, apportant régulièrement de nouvelles possibilités ainsi que des défis techniques et éthiques à considérer.

5 Bibliographie

- 1) Esser, Patrick, et al. *Scaling Rectified Flow Transformers for High-Resolution Image Synthesis*. 2024, <https://arxiv.org/abs/2403.03206>.
- 2) Louapre, David. *Comment Fonctionne ChatGPT ?* <https://scienceetonnante.com/2023/04/14/comment-fonctionne-chatgpt/>. Accessed 4 Feb. 2025.
- 3) Sabatou, Alexandre, et al. *ChatGPT, Et Après ? Bilan Et Perspectives de l'intelligence Artificielle*. 2024, <https://www.senat.fr/notice-rapport/2024/r24-170-notice.html>.
- 4) Karpathy, Andrej. *Deep Dive into LLMs Like ChatGPT*. <https://www.youtube.com/watch?v=7xTGNNLPyMI>. Accessed 25 Feb. 2025.
- 5) CNRS. *Fidle*. 2025, <https://fidle.cnrs.fr/w3/>.
- 6) CNRS, IDRIS -. *Panoram'IA*. <https://www.youtube.com/@idriscnrs>. Accessed 12 Mar. 2025.
- 7) Jamil, Umar. *You Tube*. <https://www.youtube.com/@umarjamilai>. Accessed 5 Mar. 2025.
- 8) Liu, Yiheng, et al. *Understanding LLMs : A Comprehensive Overview from Training to Inference*. 2024, <https://arxiv.org/abs/2401.02038>.
- 9) *Consommation Quotidienne Brute Régionale*. <https://odre.opendatasoft.com/explore/dataset/consommation-quotidienne-brute-regionale>. Accessed 3 Feb. 2025.
- 10) *Température Quotidienne Régionale*. <https://odre.opendatasoft.com/explore/dataset/temperature-quotidienne-regionale>. Accessed 3 Feb. 2025.
- 11) *Décomposition Biais Variance En Machine Learning*. https://mariusgareaux.github.io/blog/biais_var/. Accessed 8 Dec. 2025.
- 12) Redmon, Joseph, et al. "YOLOv3 : An Incremental Improvement." *arXiv*, 2018, <https://pjreddie.com/darknet/yolo/>.
- 13) Véry, Thibaut. *AlphaFold2 Et La Prédiction de Structures de Protéines*. <https://www.youtube.com/watch?v=YEEXMsowhxg>. Accessed 6 Feb. 2025.
- 14) *AlphaFold - Probable Disease Resistance Protein At1g58602 - Google Deep Mind, EMBL-EBI*. <https://alphafold.com/entry/Q8W3K0>. Accessed 3 Feb. 2025.
- 15) *Pause IA*. <https://pauseia.fr/dangers/economiques-et-materiels>. Accessed 8 Dec. 2025.
- 16) <https://time.com/6247678/openai-chatgpt-kenya-workers/>. Accessed 4 Feb. 2025.
- 17) *Google Turns to Nuclear to Power AI Data Centres*. <https://www.bbc.com/news/articles/c748gn94k95o>. Accessed 6 Feb. 2025.
- 18) Lefèvre, Laurent, et al. *IA Générative Sobre : Un Oxymore ?* <https://2024.jres.org/programme#modal-119>. Accessed 6 Feb. 2025.
- 19) *Epoch AI*. <https://epoch.ai>. Accessed 8 Dec. 2025.
- 20) *Fiche de Criticité - Tantale*. <https://www.mineralinfo.fr/sites/default/files/documents/2020-12/fichecriticiteantale-publique20200131.pdf>. Accessed 6 Feb. 2025.
- 21) *Du Sang Dans Nos Cellulaires*. <https://ici.radio-canada.ca/info/2019/05/coltan-republique-democratique-congo-mines-enfants/>. Accessed.
- 22) *Tom's Hardware*. <https://www.tomshardware.com/pc-components/gpus/datacenter-gpu-service-life-cannot-be-surprisingly-short-only-one-to-three-years-is-expected-according-to-unnamed-google-architect>. Accessed 6 Feb. 2025.
- 23) *NVIDIA*. <https://www.nvidia.com/en-us/sustainability/product-recycling/>. Accessed 6 Feb. 2025.
- 24) *NVIDIA Corporate Sustainability Report*. <https://images.nvidia.com/aem-dam/Solutions/documents/FY2024-NVIDIA-Corporate-Sustainability-Report.pdf>. Accessed 6 Feb. 2025.
- 25) *OpenAI Pleads That It Can't Make Money Without Using Copyrighted Materials for Free*. <https://futurism.com/the-byte/openai-copyrighted-material-parliament>. Accessed 6 Feb. 2025.
- 26) Cam, Stéphanie Le. *Intelligence Artificielle Et Métiers de La Création En Bretagne : Enjeux Et Préconisations*. <https://video.univ-rennes1.fr/permalink/v126b083de664cw01kj3/iframe/>. Accessed 8 Dec. 2025.
- 27) *Can AI Scaling Continue Through 2030 ?* <https://epoch.ai/blog/can-ai-scaling-continue-through-2030#what-constraint-is-the-most-limiting>. Accessed 8 Dec. 2025.
- 28) Touvron, Hugo, et al. *LLaMA : Open and Efficient Foundation Language Models*. 2023, <https://arxiv.org/abs/2302.13971>.
- 29) Luccioni, Sasha, et al. "Power Hungry Processing : Watts Driving the Cost of AI Deployment ?" *The 2024*

- ACM Conference on Fairness, Accountability, and Transparency, ACM, 2024, <https://doi.org/10.1145/3630106.3658542>.
- 30) *Consommation d'eau de ChatGPT*. <https://apnews.com/article/chatgpt-gpt4-iowa-ai-water-consumption-microsoft-f551fde98083d17a7e8d904f8be822c4>. Accessed 6 Feb. 2025.
 - 31) *Face à l'IA Générative, l'objection de Conscience*. <https://atecopol.hypotheses.org/12971>.
 - 32) Baydin, Atilim Gunes, et al. *Automatic Differentiation in Machine Learning : A Survey*. 2018, <https://arxiv.org/abs/1502.05767>.
 - 33) *PyTorch*. <https://pytorch.org/>. Accessed 3 Feb. 2025.
 - 34) DeepSeek-AI, et al. *DeepSeek-V3 Technical Report*. 2024, <https://arxiv.org/abs/2412.19437>.
 - 35) Devlin, Jacob, et al. *BERT : Pre-Training of Deep Bidirectional Transformers for Language Understanding*. 2019, <https://arxiv.org/abs/1810.04805>.
 - 36) *GPT Tokenizer*. <https://gpt-tokenizer.dev/>. Accessed 3 Feb. 2025.
 - 37) *Llama-3.2-1B*. <https://huggingface.co/meta-llama/Llama-3.2-1B>. Accessed 4 Feb. 2025.
 - 38) *How to Generate Text : Using Different Decoding Methods for Language Generation with Transformers*. <https://huggingface.co/blog/how-to-generate>. Accessed 7 Mar. 2025.
 - 39) *Mécanisme d'attention*. <https://ai.stackexchange.com/questions/21389/what-is-the-intuition-behind-the-attention-mechanism>. Accessed 3 Feb. 2025.
 - 40) *Wikipédia*. <https://wikipedia.fr/>. Accessed 26 Feb. 2025.
 - 41) *GitHub*. <https://github.com/>. Accessed 26 Feb. 2025.
 - 42) *ArXiv*. <https://arxiv.org/>. Accessed 26 Feb. 2025.
 - 43) *FineWeb*. <https://huggingface.co/datasets/HuggingFaceFW/fineweb>. Accessed 25 Feb. 2025.
 - 44) Tie, Guiyao, et al. *A Survey on Post-Training of Large Language Models*. 2025, <https://arxiv.org/abs/2503.06072>.
 - 45) Köpf, Andreas, et al. *OpenAssistant Conversations – Democratizing Large Language Model Alignment*. 2023, <https://arxiv.org/abs/2304.07327>.
 - 46) *Tiktokenizer*. <https://tiktokenizer.vercel.app/>. Accessed 26 Feb. 2025.
 - 47) Ouyang, Long, et al. *Training Language Models to Follow Instructions with Human Feedback*. 2022, <https://arxiv.org/abs/2203.02155>.
 - 48) Mikolov, Tomas, et al. *Efficient Estimation of Word Representations in Vector Space*. 2013, <https://arxiv.org/abs/1301.3781>.
 - 49) Pennington, Jeffrey, et al. "GloVe : Global Vectors for Word Representation." *Empirical Methods in Natural Language Processing (EMNLP)*, 2014, pp. 1532–43, <http://www.aclweb.org/anthology/D14-1162>.
 - 50) Vaswani, Ashish, et al. *Attention Is All You Need*. 2023, <https://arxiv.org/abs/1706.03762>.
 - 51) Su, Jianlin, et al. *RoFormer : Enhanced Transformer with Rotary Position Embedding*. 2023, <https://arxiv.org/abs/2104.09864>.
 - 52) Jamil, Umar. *LLaMA Explained : KV-Cache, Rotary Positional Embedding, RMS Norm, Grouped Query Attention, SwiGLU*. https://www.youtube.com/watch?v=Mn_9W1nCFLo. Accessed 27 Mar. 2025.
 - 53) *Llama Nuts and Bolts - RoPE*. <https://adalkiran.github.io/llama-nuts-and-bolts/10-ROPE-ROTARY-POSITIONAL-EMBEDDINGS/>. Accessed 7 Mar. 2025.
 - 54) *Train 400x Faster Static Embedding Models with Sentence Transformers*. <https://huggingface.co/blog/static-embeddings>. Accessed 7 Mar. 2025.
 - 55) *Model2Vec : Distill a Small Fast Model from Any Sentence Transformer*. <https://huggingface.co/blog/Pringled/model2vec>. Accessed 7 Mar. 2025.
 - 56) Groeneveld, Dirk, et al. *OLMo : Accelerating the Science of Language Models*. 2024, <https://arxiv.org/abs/2402.00838>.
 - 57) DeepSeek-AI. "DeepSeek-R1 : Incentivizing Reasoning Capability in LLMs via Reinforcement Learning." *arXiv Preprint arXiv :2501.12948*, 2024, <https://arxiv.org/abs/2501.12948>.
 - 58) Grattafiori, Aaron, et al. *The Llama 3 Herd of Models*. 2024, <https://arxiv.org/abs/2407.21783>.
 - 59) *Opening up ChatGPT* *Opening up ChatGPT : Tracking Openness of Instruction-Tuned LLMs*. <https://opening-up-chatgpt.github.io/>. Accessed 18 Mar. 2025.
 - 60) *Open Source AI Index*. <https://osai-index.eu/the-index>. Accessed 21 Mar. 2025.
 - 61) OpenAI. *GPT-4 Technical Report*. 2023, <https://arxiv.org/abs/2303.08774>.
 - 62) Anthropic. *Claude*. 2024, <https://www.anthropic.com/news/claude-3-family>.
 - 63) Touvron, Hugo, et al. *LLaMA : Open and Efficient Foundation Language Models*. 2023, <https://arxiv.org/abs/2302.13013>.

- g/abs/2302.13971.
- 64) Team, Gemini, et al. *Gemini : A Family of Highly Capable Multimodal Models*. 2024, <https://arxiv.org/abs/2312.11805>.
 - 65) Chen, Stanley F., et al. "An Empirical Study of Smoothing Techniques for Language Modeling." *Computer Speech & Language*, vol. 13, no. 4, 1999, pp. 359–94, <https://aclanthology.org/A98-1014.pdf>.
 - 66) Shannon, Claude E. "A Mathematical Theory of Communication." *Bell System Technical Journal*, vol. 27, no. 3, 1948, pp. 379–423, <https://people.math.harvard.edu/~ctm/home/text/others/shannon/entropy/entropy.pdf>.
 - 67) Brown, Peter F., et al. "Class-Based n-Gram Models of Natural Language." *Proceedings of the 30th Annual Meeting of the Association for Computational Linguistics*, 1992, pp. 467–73, <https://aclanthology.org/P92-1058.pdf>.
 - 68) Elman, Jeffrey L. "Finding Structure in Time." *Cognitive Science*, vol. 14, no. 2, 1990, pp. 179–211, [https://doi.org/10.1016/0364-0213\(90\)90002-E](https://doi.org/10.1016/0364-0213(90)90002-E).
 - 69) Hochreiter, Sepp, et al. "Long Short-Term Memory." *Neural Computation*, vol. 9, no. 8, 1997, pp. 1735–80, <https://www.bioinf.jku.at/publications/older/2604.pdf>.
 - 70) Cho, KyungHyun, et al. *Learning Phrase Representations Using RNN Encoder-Decoder for Statistical Machine Translation*. 2014, <https://arxiv.org/abs/1406.1078>.
 - 71) Bengio, Yoshua, et al. "A Neural Probabilistic Language Model." *Journal of Machine Learning Research*, vol. 3, 2003, pp. 1137–55, <https://www.jmlr.org/papers/volume3/bengio03a/bengio03a.pdf>.
 - 72) Brants, Thorsten, et al. "Large Language Models in Machine Translation." *Proceedings of the 2007 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning*, 2007, pp. 858–67, <https://aclanthology.org/D07-1090.pdf>.
 - 73) Mikolov, Tomáš, et al. "Recurrent Neural Network Based Language Model." *Proceedings of the 11th Annual Conference of the International Speech Communication Association (INTERSPEECH 2010)*, 2010, pp. 1045–48, https://www.isca-speech.org/archive/interspeech_2010/mikolov10_interspeech.html.
 - 74) Bahdanau, Dzmitry, et al. *Neural Machine Translation by Jointly Learning to Align and Translate*. 2014, <https://arxiv.org/abs/1409.0473>.
 - 75) Radford, Alec, et al. *Improving Language Understanding by Generative Pre-Training*. OpenAI, 2018, https://cdn.openai.com/research-covers/language-unsupervised/language_understanding_paper.pdf.
 - 76) Radford, Alec, et al. *Language Models Are Unsupervised Multitask Learners*. OpenAI, 2019, https://cdn.openai.com/better-language-models/language_models_are_unsupervised_multitask_learners.pdf.
 - 77) Yang, Zhilin, et al. *XLNet : Generalized Autoregressive Pretraining for Language Understanding*. 2019, <https://arxiv.org/abs/1906.08237>.
 - 78) Liu, Yinhan, et al. *RoBERTa : A Robustly Optimized BERT Pretraining Approach*. 2019, <https://arxiv.org/abs/1907.11692>.
 - 79) Brown, Tom, et al. "Language Models Are Few-Shot Learners." *Advances in Neural Information Processing Systems*, vol. 33, 2020, pp. 1877–901, https://proceedings.neurips.cc/paper_files/paper/2020/file/1457c0d6bfc4967418bf8ac142f64a-Paper.pdf.
 - 80) Raffel, Colin, et al. "Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer." *Journal of Machine Learning Research*, vol. 21, no. 1, 2020, pp. 5485–551, <https://jmlr.org/papers/volume21/20-074/20-074.pdf>.
 - 81) Shoeybi, Mohammad, et al. *Megatron-LM : Training Multi-Billion Parameter Language Models Using Model Parallelism*. 2019, <https://arxiv.org/abs/1909.08053>.
 - 82) Stability AI. *StableLM-3B-4E1T*. 2024, <https://stability.wandb.io/stability-llm/stable-lm/reports/StableLM-3B-4E1T--VmllDzoyMjU4>.
 - 83) Alayrac, Jean-Baptiste, et al. *Flamingo : A Visual Language Model for Few-Shot Learning*. 2022, <https://arxiv.org/abs/2204.14198>.
 - 84) Driess, Danny, et al. *PaLM-e : An Embodied Multimodal Language Model*. 2023, <https://arxiv.org/abs/2303.03378>.
 - 85) Liu, Haotian, et al. *Visual Instruction Tuning*. 2023, <https://arxiv.org/abs/2304.08485>.
 - 86) Chen, Mark, et al. *Evaluating Large Language Models Trained on Code*. 2021, <https://arxiv.org/abs/2107.03374>.
 - 87) Rozière, Baptiste, et al. *Code Llama : Open Foundation Models for Code*. Meta AI, 2023, <https://ai.meta.com/blog/code-llama-large-language-model-coding/>.

- 88) Singhal, Karan, et al. *Towards Expert-Level Medical Question Answering with Large Language Models*. 2023, <https://arxiv.org/abs/2305.09617>.
- 89) Lee, Jinhyuk, et al. "BioBERT : A Pre-Trained Biomedical Language Representation Model for Biomedical Text Mining." *Bioinformatics*, vol. 36, no. 4, 2020, pp. 1234–40, <https://doi.org/10.1093/bioinformatics/btz682>.
- 90) Chalkidis, Ilias, et al. "LEGAL-BERT : The Muppets Straight Out of Law School." *Findings of the Association for Computational Linguistics : EMNLP 2020*, Association for Computational Linguistics, 2020, pp. 2898–904, <https://doi.org/10.18653/v1/2020.findings-emnlp.261>.
- 91) *Ollama - Run Llama 2, Mistral, and Other Models Locally*. <https://ollama.ai/>. Accessed 28 Mar. 2025.
- 92) *LM Studio - Discover, Download, and Run Local LLMs*. <https://lmstudio.ai/>. Accessed 28 Mar. 2025.
- 93) Gerganov, Georgi. *Llama.cpp - Port of Facebook's LLaMA Model in c/c++*. <https://github.com/ggerganov/llama.cpp>. Accessed 28 Mar. 2025.
- 94) *Replicate - Run Open-Source Models with an API*. <https://replicate.com/>. Accessed 28 Mar. 2025.