

Fine-tuning des LLM

Formation continue – 1 jour

Fine-tuning : ajustement / réglage fin

Idée : Continuer l'entraînement d'un grand modèle de langage, par exemple pour effectuer une tâche spécifique :

- classification de texte,
- spécification de vocabulaire dans un domaine,
- synthèse de texte

Modèles ajustés (« *fine-tunés* »)

Transforme du texte en
requêtes SQL (base de
données)



Modèle de
fondation



agent-conversationnel



Agent conversationnel
spécialisé dans un
domaine (connaît un
vocabulaire spécifique)



Agent conversationnel
très bon en synthèse de
texte



Modèle de classification
de phrases



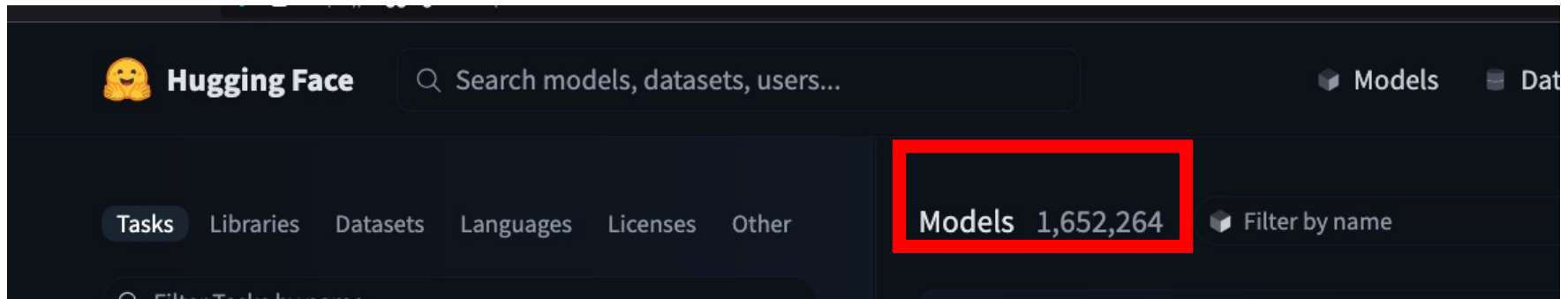
Fine-tuning des LLM

En pratique, plusieurs familles de modèles de fondation (avec plusieurs types de modèles : textes, audio, vision, multimodaux, ...) :

- Llama,
- BERT,
- GPT,
- Qwen,
- DeepSeek,
- BLOOM

....

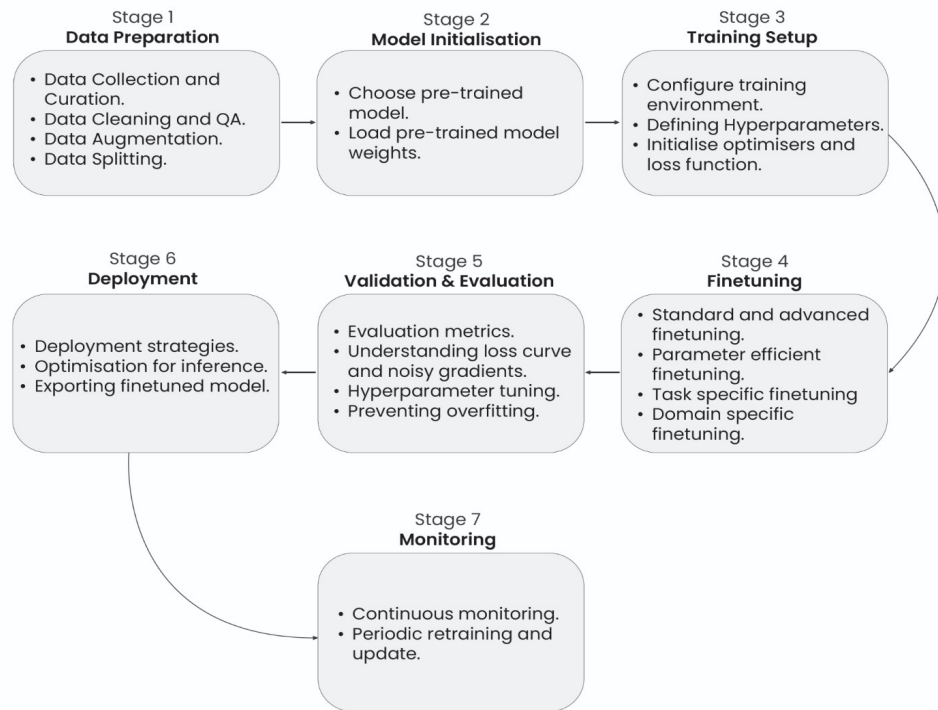
Et beaucoup de *fine-tuning* de ces modèles (<https://huggingface.co/models>):



Quelques ingrédients nécessaires :

- des données (**structurées!**),
- des moyens de calculs (local ou distant),
- un modèle pré-entraîné,
- les outils permettant d'entraîner le modèle sur les données (bibliothèque Python, accès haut niveau, gestion des sessions d'entraînements, ...),
- une méthode d'entraînement (mise à jour de tous les paramètres, seulement certains, LoRA, ...)

Fine-tuning des LLM



Source : <https://arxiv.org/pdf/2408.13296v1>

les données :

- dépend de l'objectif de l'entraînement :
 - résumé de textes : des listes de textes et leurs résumés,
 - classification de sentiments : des phrases et des scores,
 - agent conversationnel : des exemples de conversations
- jeux de données publics (*potentiellement sous licence*) :
 - [HuggingFace Datasets](#)
- données propres à votre domaine :
 - bases de données **bien indexées**

Dans tous les cas, les données peuvent être totalement ou partiellement générées par des LLMs

choisir l'infrastructure:

- **Local** : possible selon la / les machines à disposition (carte graphique / puces Apple Silicon préférables). Le fine-tuning peut ne pas demander beaucoup de ressources (~quelques heures d'entraînement sur un PC perso)
- **Distant** : (source <https://towardsai.net>)
 - sur des plateformes proposant des accès distants à des machines :
 - classiques : Google Cloud Platform, Microsoft Azure, ..
 - **plus locales** : IDRIS, Onyxia, ...
 - *via une API de fine-tuning* (ex : OpenAI)

choisir les outils :

- Avec une / plusieurs bibliothèques Python :
 - [HuggingFace Trainer](#), et tous ses dérivés (beaucoup de modèles différents supportés)
 - Autres bibliothèques plus spécifiques, mais plus simples d'accès :
 - [Llama CookBook](#),
 - [Unsloth](#),
 - ...
- *Via* un accès haut niveau (ex : API de fine-tuning OpenAI) :
 - Pas/ peu la main sur l'entraînement
 - Diffusion des modèles bloquée par le propriétaire de la plateforme



Fine-tuning des LLM

choisir les outils :

- plateforme d'entraînement (gestion des « runs », gestionnaire de workflow, analyse des performances du modèle entraîné)
- quelques exemples (local ou distant) :
 - AirFlow, Prefect, MLFlow, MetaFlow, ...

choisir un modèle :

- le *fine-tuning* passe par la **modification des paramètres** du modèle. Impossible pour les modèles à tendance fermés, accessibles uniquement *via API* – à moins d'avoir accès à une API pour le *fine-tuning*
 - privilégier les modèles à **tendance ouverte**
- le choix du modèle peut être conditionné par les outils utilisés pour le *fine-tuning*.
 - ex : **Llama CookBook** n'accepte que les modèles Llama
- taille du modèle : incidence sur le temps d'entraînement (+de paramètres => +long).

choisir des méthodes de fine-tuning (complémentaires!):

- Supervised Fine-Tuning,
- LoRA,
- PEFT,
- Reinforcement Learning with Human Feedback (RLHF),
- DPO,
- QLoRA,
- ...

→ différent selon les cas d'usages, et sont plus ou moins complexes à mettre en place

Fine-tuning des LLM

Exemple / expérimentation : Intégrer des acronymes et leurs définitions dans un agent conversationnel

- 0 – Avoir une liste **structurée** d'acronymes et leurs définitions
- 1 – Générer des conversations (avec un LLM par exemple) sur ces acronymes
- 2 – Continuer l'entraînement d'un agent conversationnel sur ces données
 - choisir un modèle
 - choix de l'infrastructure d'entraînement (local, distant)
 - méthode et outils pour le *fine-tuning*
- 3 – Tester le modèle
- 4 – Partager, utiliser le modèle *fine-tuné*

0 – Avoir une liste **structurée** d'acronymes et leurs définitions

```

73  },
74  {
75    "acronym": "PER",
76    "definition": "Purée et Épice de la Réussite"
77  },
78  {
79    "acronym": "FUN",
80    "definition": "Fusion Universelle de la Nourriture"
81  },
82  {
83    "acronym": "TEC",
84    "definition": "Techniques Élémentaires et Culinaires"
85  },
86  {
87    "acronym": "SAL",
88    "definition": "Savoir-Faire et Aromes de la Légèreté"
89  },
90  {
91    "acronym": "CAMP",
92    "definition": "Création Artisanale et Maîtrisée de la Pâtisserie"
93  },

```

```

1423
1424  <data skos:Concept="http://www.my.com/#institute" uri="http://www.my.com/#kisti">
1425    <skos:prefLabel>kisti</skos:prefLabel>
1426    <skos:altLabel>Korean Institute for Science and Technology Information</skos:altLabel>
1427    <data uri="http://www.my.com/#korea"/>
1428    <skos:related rdf:resource="http://www.my.com/#research"/>
1429    <skos:definition>KISTI is a government-funded research institute designed to maximize the efficiency of science and t
1430    <link skos:Concept="http://www.my.com/#url" value="http://en.kisti.re.kr"/>
1431  </data>
1432
1433  <data skos:Concept="http://www.my.com/#laboratory" uri="http://www.my.com/#BerkeleyLab">
1434    <skos:prefLabel>Lawrence Berkeley National Laboratory</skos:prefLabel>
1435    <data uri="http://www.my.com/#USA"/>
1436    <data uri="http://www.my.com/#AnshuDubey"/>
1437    <skos:related rdf:resource="http://www.my.com/#research"/>
1438    <skos:definition>In the world of science, Lawrence Berkeley National Laboratory (Berkeley Lab) is synonymous with "e
1439    <link skos:Concept="http://www.my.com/#url" value="http://www.lbl.gov"/>@
1440  </data>
1441  <data skos:Concept="http://www.my.com/#supercomputer" uri="http://www.my.com/#Yellowstone">
1442    <data uri="http://www.my.com/#USA"/>
1443    <data uri="http://www.my.com/#NCAR"/>
1444    <skos:definition>Yellowstone is NCAR's 1.5-petaflops high-performance IBM iDataPlex cluster, which features 72,576 Ib
1445    <link skos:Concept="http://www.my.com/#url" value="https://www2.cisl.ucar.edu/resources/yellowstone"/>
1446  </data>
1447

```

1 – Générer des conversations (avec un LLM par exemple) sur ces acronymes

gemma3:27b v +

Create 10 fictive conversations IN FRENCH between an user and an assistant. Those conversations must contains 1 question and 1 answer. Each question must be an user asking for the definition of the acronym TEC; and each answer must contain the definition : 'Techniques Élémentaires et Culinaires'; or a more verbose definition : Techniques Élémentaires et Culinaires. All the answer must be somehow diverse. Each conversation will be formatted in a json list, where each element is itself a list of the form : [{ 'role': 'user', 'content': THE QUESTION }, { 'role': 'assistant', 'content': THE ANSWER }] Keep it short. The answer must be the raw json; no fioritures.

```
{
  "role": "user",
  "content": "Bonjour, pouvez-vous me dire ce que signifie l'acronyme TEC ?"
},
{
  "role": "assistant",
  "content": "TEC signifie Techniques \u00c9l\u00e9mentaires et Culinaire."
}
```



2 - Entraînement

Mini-cours sur l'entraînement et le fonctionnement des réseaux de neurones

<https://www.figma.com/design/bdHDd6JYnM5Dfsoi8lhBwP/Untitled?m=auto&t=LQAal7R3>



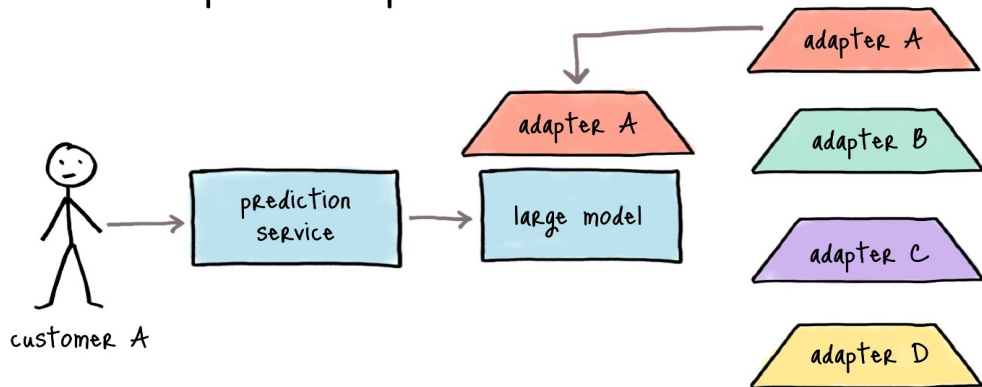
Mini-cours sur l'entraînement et le fonctionnement des réseaux de neurones

Résumé du vocabulaire :

- learning-rate : « pas » du gradient, mesure la vitesse d'apprentissage
- n_epochs : nombre de fois où le modèle a « vu » toute les données d'entraînement
- mini-batch-size (per_device_train_batch_size dans « transformers »): nombre d'éléments du jeu de données utilisés dans 1 mise à jour des poids du modèle

Zoom sur LoRA (Low-Rank Adaptation) :

LoRA adapters - personalized models



Source :
<https://huggingface.co/learn/lm-course/en/chapter11/4>

Idée : On **rajoute** des **paramètres**, que l'on entraîne avec les données. Le nombre de paramètres rajoutés est bien **plus faible** (par ex. 5%) que le nombre de paramètres du modèle.

- entraînement plus rapide ; et donne des résultats *presque* aussi bons,
- sépare le modèle de base (**lourd**) de la partie *fine-tuning* (**léger**),
- permet d'utiliser plusieurs *fine-tuning* différents sur un même modèle.

Couche d'un modèle pré-entraîné :

$h = W.x$ (W : matrice avec d lignes et d colonnes)

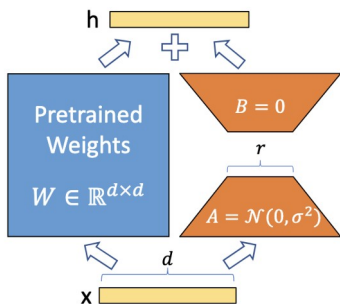
Fine-Tuning « classique » :

$h = W'.x$ (on modifie W en W' pendant le fine-tuning)

$h = W.x + \Delta W.x$ ($\Delta W = W' - W$)
 $\rightarrow \Delta W$ et W ont la même taille (d lignes, d colonnes ! $\rightarrow d^2$ paramètres mis à jours. Si $d = 3072$, cela fait : 9 437 184)

Fine-tuning LoRA :

$h = W.x + B.A.x$



Source:
<https://arxiv.org/pdf/2106.0968>

- $\rightarrow A$ et B ont des dimensions plus petites que W
- \rightarrow On choisit A de taille (d, r) et B de taille (r, d) .
 (A a d lignes et r colonnes, B a r lignes et d colonnes)
 \rightarrow les coefficients de A et B seront les seuls modifiés pendant l'entraînement
- \rightarrow Typiquement, on fixe $r = 16$.
- $\rightarrow 2.d.r$ paramètres : (si $d = 3072$, cela fait : 98 304)



Exemple / expérimentation : Mémorisation d'acronymes

2 – Continuer l'entraînement d'un agent conversationnel sur ces données

choisir un modèle



Llama 3.2-1B
Instruct

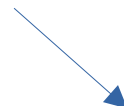
- petit modèle, donc rapide à entraîner
- à tendance ouverte, possible de le *fine-tuner* localement
- ! licence de diffusion restrictive

choisir une infrastructure



- local (Mac M3) :
pour expérimenter
- distant (GCP / Onyxia) :
entraînement plus important (**mais pas nécessaire ici – facteur x3 temps**)

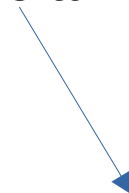
outils d'entraînement



**Hugging Face
Trainer :**

- gestion fine des hyper-paramètres,
 - bien documenté
 - prise en main très simple
 - gratuit + open-source
- Metaflow, MLFlow
(gestionnaires de workflow)**

méthode de *fine-tuning*



- LoRA : **léger et modulaire**

3 – Tester le modèle

On veut tester le modèle sur sa capacité à **mémoriser des définitions**.

On pose au modèle les questions du jeu de test, et on évalue ensuite la qualité de ses réponses.

2 exemples pour évaluer la qualité des réponses :

- **modèles d'« embeddings »** : petits modèles, qui sont entraîné à prédire le taux de similarité entre deux séquences de mots

- **« LLM as a judge »** : on demande à un LLM tiers d'évaluer la qualité des réponses du LLM « fine-tuné ». Par exemple : « Réponds 1 si les deux réponses sont proches, et 0 sinon ».



3 – Tester le modèle

Selon les résultats du modèles, recommencer l'entraînement avec des paramètres différents .. !

4 – Partager et utiliser le modèle

Quelques exemples :

- publier sur Hugging Face Hub
 - type d'accès à décider : accès ouvert, accès restreint
 - licence à décider
- partager le modèle brut :
 - poids : .safetensors (~pkl) ou autres formats quantisés
 - directement les checkpoints de l'entraînement pour l'expérimentation
- dans le cas d'un modèle fine-tuné avec la méthode LoRA on peut se contenter de partager la surcouche (+ légère que le modèle entier)



Exemple / expérimentation : Mémorisation d'acronymes

Pour expérimenter :

→ code disponible sur github :

https://github.com/mariusgarenaux/fine_tuning_acronym

FOLLOW US!

