

Fine-Tuning LoRA/QLoRA : Guide P

29 April
2026Mis à jour le 29 April
202648 min de
lecture

Guide complet fine-tuning LoRA/QLoRA : PEFT, Unsloth, Axolotl, dataset A
vLLM/TGI. Comparatif RAG vs fine-tuning.

Le fine-tuning des grands modèles de langage est devenu accessible grâce à des méthodes qui réduisent drastiquement les ressources nécessaires. LoRA (Low-Rank Adaptation) et les méthodes PEFT (Parameter-Efficient Fine-Tuning) permettent de spécialiser des modèles sur un seul GPU consommateur. Alors qu'un fine-tuning complet de Llama 3 70B nécessiterait des milliers de dollars de calcul, LoRA permet d'obtenir des résultats comparables en ne modifiant que la VRAM et quelques dizaines de dollars. Ce guide exhaustif couvre les fondements de LoRA, la préparation des datasets (formats Alpaca et ShareGPT), les outils d'entraînement (vLLM, TGI, OpenAI harness), le déploiement (vLLM, TGI), la comparaison avec le RAG et le prompt engineering. Pour les praticiens de l'IA, les data scientists et les ingénieurs ML, ce guide est une base fondamentale pour personnaliser les LLM à des cas d'usage spécifiques sans expertise avancée.

À RETENIR

A retenir : LoRA ne modifie que 0,1 à 1 % des paramètres du modèle, réduisant ainsi le coût de calcul. Combiner la quantification 4 bits et LoRA pour fine-tuner des modèles 70B sur une tâche spécifique est une approche de fine-tuning complète pour la plupart des cas d'usage, à une fraction du coût.

Pourquoi fine-tuner un LLM ?

Le fine-tuning consiste à poursuivre l'entraînement d'un modèle pré-entraîné sur un style ou une tâche particulière. Avant d'explorer comment, il est essentiel de comprendre pourquoi le fine-tuning est la bonne approche pour votre cas d'usage.

Quand le fine-tuning est-il nécessaire ?

Le fine-tuning est justifié dans plusieurs scénarios. Pour l'adaptation de style ou d'ajout de connaissances d'une manière spécifique (jargon technique, style juridique, ton de marque) que le modèle générique ne peut pas fournir. Pour l'apprentissage de formats de sortie complexes, quand le modèle doit systématiquement produire des résultats précis (JSON structure, XML spécifique, format de rapport particulier). Pour la spécialisation dans un domaine (médical, juridique) et maîtriser un vocabulaire et des concepts spécifiques à un domaine (médical, juridique) que le modèle générique ne maîtrise pas. Pour la réduction de latence et de coût, quand vous recherchez les performances d'un modèle plus grand sur votre tâche spécifique. Et pour l'intégration de connaissances, si vous voulez que le modèle "connaisse" des informations spécifiques à votre organisation.

Fine-tuning vs RAG vs Prompt Engineering

Le choix entre fine-tuning, RAG et prompt engineering dépend du problème à résoudre. L'approche à essayer : il est gratuit (pas de coût d'entraînement), instantané et réversible.
