

vLLM vs Ollama vs TGI vs SGLang — bench



16 mai
2026



Mis à jour le 17 mai
2026



16 min de
lecture



3133
mots

Benchmark complet des serveurs LLM en 2026 : vLLM, Ollama, TGI, SGLang, consommation GPU, facilité de déploiement. Quel serveur choisir selon vos besoins ?

À RETENIR

A retenir -- Benchmark serveurs LLM 2026

En 2026, quatre frameworks dominent le **servicing LLM en production** : vLLM (performance), Ollama (simplicité de déploiement sur GPU grand public), TGI (robustesse et consommation GPU minimale pour les workloads interactifs). Le choix dépend de trois facteurs : la disponibilité des ressources, la priorité entre throughput et latence, et la facilité de déploiement. Pour 90% des déploiements en production avec plus de 10 requêtes concurrentes, vLLM est l'option optimale.

Le choix du **framework de servicing LLM** est une décision d'infrastructure critique qui impacte directement les performances et la complexité opérationnelle d'un déploiement IA en production. En 2026, l'écosystème est dominé par quatre acteurs principaux -- chacun

Réponse sous 24h

principaux -- vLLM, Ollama, TGI (Text Generation Inference), SGLang

Devis
gratuit



d'usage. Ce benchmark compare ces quatre solutions sur les metriques les plus importantes : le throughput (tokens/seconde), latence P50 et P99, efficacite memoire GPU, facilite d'integration avec les modeles open source majeurs (Llama 3.3, Mistral, Qwen 2.5), et niveau de maturite de la configuration representative (GPU NVIDIA A100 40GB SXM) avec le modele Llama 3.3 et en float16 pour vLLM, TGI et SGLang.

vLLM -- le champion du throughput haute performance

vLLM (docs.vllm.ai) est un framework de serving LLM developpe par UC Berkeley avec support commercial. Son innovation cle est le **PagedAttention** : une gestion innovante des systemes d'exploitation, qui permet d'utiliser la memoire GPU de facon beaucoup plus efficace en dynamiquement des "pages" de KV-cache selon les besoins de chaque sequence.

Les avantages de vLLM en production :

Throughput maximal : vLLM est generalement 2 a 4x plus rapide que les implémentations open source, plus rapide que TGI sur des charges elevees avec de nombreuses requetes concurrentes.

Continuous batching : contrairement au batching statique, vLLM traite en continu les requetes entrant, tant que le courant soit termine, reduisant considerablement la latence percue sous charge.

API OpenAI compatible : drop-in replacement pour l'API OpenAI, facilitant la migration des applications existantes.

Support multi-GPU : tensor parallelism et pipeline parallelism pour les modeles de grande taille.

```
# Installation et demarrage vLLM
```

```
pip install vllm
```

Un projet cybersécurité ?

Réponse sous 24h

```
# Serveur OpenAI-compatible sur GPU A100
```

Devis
gratuit →

Réponse sous 24h

Devis
gratuit →