



# vLLM : Moteur d'Inférence LLM Haute Performances 2026



10 mai 2026



Mis à jour le 17 mai 2026



23 min de lecture



4748 mots



vLLM est un moteur open-source d'inférence et de service pour LLM, écrit en Python et CUDA, conçu pour offrir un débit maximal et une latence prévisible sur GPU et accélérateurs spécialisés. Né en 2023 au Sky Computing Lab de UC Berkeley sous l'impulsion de Woosuk Kwon, Zhuohan Li, Ion Stoica et Hao Zhang, vLLM cumule en mai 2026 plus de 52 000 étoiles GitHub, 1 100 contributeurs et est une partie de la PyTorch Foundation. Cette page entity-first détaille PagedAttention, le continuous batching, l'architecture worker/scheduler/executor, les 250+ architectures supportées (Llama 4, Mistral, Mixtral, Qwen 3, DeepSeek V3, Gemma 3, GLM-4.5), les backends CUDA/ROCm/CPU/TPU/Neuron/Gaudi, les formats FP8/AWQ/GPTQ/NVFP4, l'API OpenAI-compatible, le speculative decoding, le disaggregated prefill, le prefix caching, le multi-LoRA serving, la vLLM PaaS, Stack Helm Kubernetes, le monitoring Prometheus et les benchmarks face à TensorRT-LLM, llama.cpp et SGLang.

Un projet cybersécurité ?

Réponse sous 24h

Devis gratuit



**vLLM** est un **moteur d'inférence et de service** pour *large language models* (LLM) source, écrit en Python et CUDA, conçu pour offrir un **débit (throughput) maximal** et une **latence prévisible** sur GPU et accélérateurs spécialisés. Initialement développé en interne au sein du **Sky Computing Lab** de l'**Université de Californie à Berkeley** par **Woosuk Liu, Zhuohan Li** et leurs co-auteurs sous la direction des professeurs **Ion Stoica** et **Hadi Pajouh**, vLLM a été rendu public le **20 juin 2023** et a depuis été adopté comme *backend* de référence par une grande partie de l'industrie : Anyscale, Databricks, [Anthropic](#) (pour des charges internes), AWS Bedrock (sur certaines familles de modèles), [Cloudflare Workers](#), Lambda Labs, Together AI, RunPod, Mistral La Plateforme et l'écosystème Red Hat OpenShift AI. La version **v0.10.x** de mai 2026 cumule plus de **52 000 étoiles GitHub**, dépassant les autres moteurs de référence, et fait partie des projets phares de la *PyTorch Foundation*, où elle a été ajoutée en septembre 2025. Son innovation algorithmique principale, **PagedAttention**, transforme la pagination mémoire des systèmes d'exploitation au *KV cache* des transformers et offre un facteur 2 à 24 sur le débit par rapport à des moteurs naïfs (HuggingFace TGI vs vLLM, FasterTransformer historique). vLLM expose nativement une **API HTTP compatible OpenAI** (*Chat Completions, Completions, Embeddings, Tools/Function calling*), supporte plusieurs **architectures** (Llama 2/3/4, Mistral, Mixtral, Qwen 2/2.5/3, DeepSeek V2/V3/R1, Phi-3, Gemma 2/3, GLM-4, Yi, Command-R, MiniCPM, Aya, Granite, InternLM, OLMo, StarCoder, Falcon Mamba), accepte les principaux **formats de quantization** (FP16, BF16, INT8, GPTQ, FP8 E4M3/E5M2, GGUF en lecture partielle, SqueezeLLM), tourne sur **CUDA** (NVIDIA Hopper/Blackwell/Ada/Ampere), **ROCm** (AMD MI250/MI300X/MI325X), **CPU** (Intel Xeon AVX-512, ARM Neoverse), **TPU** (Google v5e/v5p/Trillium) et **AWS Neuron** (Trainium, Inferentia2). Cette page entity-first détaille l'histoire, l'algorithme PagedAttention, le *batching*, l'architecture interne (worker, scheduler, executor), la matrice de comparaison des modèles/backends/formats, l'API [OpenAI](#)-compatible, les techniques avancées (SpecDecoding, Disaggregated Prefill, Prefix Caching, multi-LoRA serving), le déploiement sur Kubernetes via la *vLLM Production Stack*, le monitoring Prometheus, les benchmarks Ollama, TensorRT-LLM, llama.cpp et SGLang, ainsi que les considérations de sécurité.

Réponse sous 24h

Devis  
gratuit



---

---

Réponse sous 24h

Devis  
gratuit

