




# FlashQuant : Compression KV-Cache 7.5x O

 10 mai  
2026Mis à jour le 17 mai  
202613 min de  
lecture270  
mot

FlashQuant est un projet open source C++17/CUDA qui compresse le KV-cache à un ratio 7.5x via quantization INT4 et codebooks par tête. Intégration native vLLM (<1%), VRAM divisée par 7.5 sur Llama 3.1 70B (40 Go → 5.3 Go).

FlashQuant est un projet open source écrit en C++17 et CUDA qui s'attaque au goulot de

critique des LLM modernes en inférence : la consommation mémoire du KV-cache. Par exemple, comme Llama 3.1 70B servi avec un contexte de 32K tokens, le KV-cache occupe la mémoire disponible, dépassant régulièrement les 40 Go par requête longue. FlashQuant introduit une compression extrême combinant quantization INT4 par groupe, dictionnaires de codebooks et CUDA dédiés, atteignant un ratio de compression moyen de 7.5x avec une dégradation de 1% sur MMLU et HumanEval. Conçu comme un plugin pour vLLM 0.6.x+, FlashQuant permet de servir des contextes longs sur GPU consumer (RTX 4090) où auparavant le modèle ne tenait précédemment que sur des H100 80 Go. Cette page projet documente les benchmarks, l'installation et la roadmap. Le code source est disponible sur [github.com/ayinedjimi/FlashQuant](https://github.com/ayinedjimi/FlashQuant).

Devis  
gratuit



## A retenir

**Ratio 7.5x** de compression KV-cache (FP16 → INT4 + codebooks) sur Llama

**VRAM divisée par 7.5** : 40 Go → 5.3 Go pour 32K tokens de contexte.

**Perplexité préservée** : dégradation inférieure à 1 % sur MMLU et HumanEval

**Latence overhead** inférieur à 5 % grâce à la déquantification on-the-fly fusionnée

**Plugin vLLM** compatible PagedAttention, sm\_80+ requis (Ampere et supérieur)

## Le problème : le KV-cache devore la VRAM

Lorsqu'un LLM autoregressif génère du texte, chaque nouveau token doit accéder à la valeur de tous les tokens précédents. Pour éviter de recalculer ces tenseurs à chaque pas, le modèle stocke un KV-cache. Sa taille croît linéairement avec la longueur de contexte, le nombre de têtes d'attention et la dimension par tête. Pour Llama 3.1 70B en FP16, le KV-cache pour 32K tokens, soit plus de 40 Go pour un contexte de 32K. Sur des modèles MoE ou des modèles à très grande échelle, cela explose au-delà des 100 Go par séquence active. Résultat : même un H100 80 Go ne peut pas fonctionner en session simultanément, et les GPU consommateurs sont totalement exclus du jeu. Cette limitation empêche le batching, augmente le coût par token et empêche le déploiement de RAG long-contexte à grande échelle et abordable.

Concrètement, le calcul est implacable : pour Llama 3.1 70B, 32K tokens, 8 K  
Attention), dimension 128 par tête, en FP16 (2 octets) consomme 80

Devis  
gratuit



---

Réponse sous 24h

Devis  
gratuit →