

# Quantization LLM — GPTQ, AWQ, EXL2,

16 mai  
2026Mis à jour le 17 mai  
202615 min de  
lecture3261  
mots15  
vue

Comparez les méthodes de quantization LLM en 2026 : GPTQ, AWQ, EXL2, VRAM requise, la vitesse d'inférence. Guide technique pour choisir la quan

## À RETENIR

### A retenir -- Quantization LLM 2026

La **quantization LLM** permet de réduire de 50 à 75% la VRAM nécessaire pour la précision des poids de float16 (16 bits) à int4 ou int8. En 2026, quatre formes sont populaires : **GGUF** (CPU friendly), **GPTQ** (GPU only, haute qualité), **AWQ** (meilleure préservation de la précision et du débit GPU). Pour un déploiement on-premise type ETI, AWQ Q4 ou GGUF Q4 sont des compromis : 70B modèle sur un GPU A100 40GB avec moins de 2% de dégradation.

La **quantization des LLM** est la technique la plus puissante pour démocratiser l'accès à l'IA sans investir dans des infrastructures GPU prohibitives. Sans quantization, Llama 3.1 nécessite 2x A100 40GB ou 2x H100 80GB). Avec quantization int4, le même modèle tourne sur

en consumer grade. Cette reduction drastique des besoins en VRAM ouvre des possibilités qui étaient inatteignables économiquement pour la majorité des organisations. En 2024, le marché s'est consolidé autour de quatre formats principaux -- GPTQ, AWQ, EXL2 et GGUF -- qui sont distinctes en termes de qualité préservée, de vitesse d'inférence, de compatibilité et de coût. Ce guide technique compare ces quatre approches pour aider les équipes MLOps à choisir la meilleure option dans leur contexte.

---

## Principes de la quantization -- réduire la précision des poids LLM

---

La **quantization** des réseaux de neurones réduit la précision numérique des paramètres, ce qui permet de réduire la consommation mémoire et accélérer les calculs. Pour les LLM, les poids sont typiquement représentés en float32 (par valeur) ou bfloat16. La quantization int8 (8 bits, 1 octet) divise la VRAM par 2, et la quantization int4 (4 bits, 0.5 octet) divise par 4.

Le challenge de la quantization est de minimiser la dégradation de qualité résultant de ces réductions de précision. Les approches naïves (arrondir chaque poids au int4 le plus proche) dégradent significativement la performance sur des tâches complexes. Les méthodes avancées (GPTQ, AWQ) utilisent des algorithmes d'optimisation pour préserver la précision en ajustant les poids adjacents ou en préservant les poids les plus importants.

**Q8\_0 / int8** : dégradation quasi-nulle (<0.5% sur les benchmarks standards), réduction VRAM par 2. Utile lorsque la VRAM est la contrainte et que la qualité est prioritaire.

**Q4\_K\_M / int4** : dégradation faible (1-2% sur les benchmarks standards), réduction VRAM par 4. Meilleur compromis qualité/VRAM pour la plupart des déploiements.

**Q2\_K / int2** : dégradation significative (5-15%), réduction VRAM par 8. Utile uniquement pour des déploiements sur GPU consumer sans suffisamment de VRAM.

---