

Llama 4, Mistral Large, Gemma 3 : Comparatif LLM Open Source

Catégorie : Intelligence Artificielle | Lecture : 18 min | Publié le : 13/02/2026 | Auteur : Ayi NEDJIMI

Comparatif détaillé des LLM open source 2026 : Llama 4, Mistral Large, Gemma 3, Qwen 2.5, DeepSeek V3. Benchmarks, coûts et guide de choix. Guide.

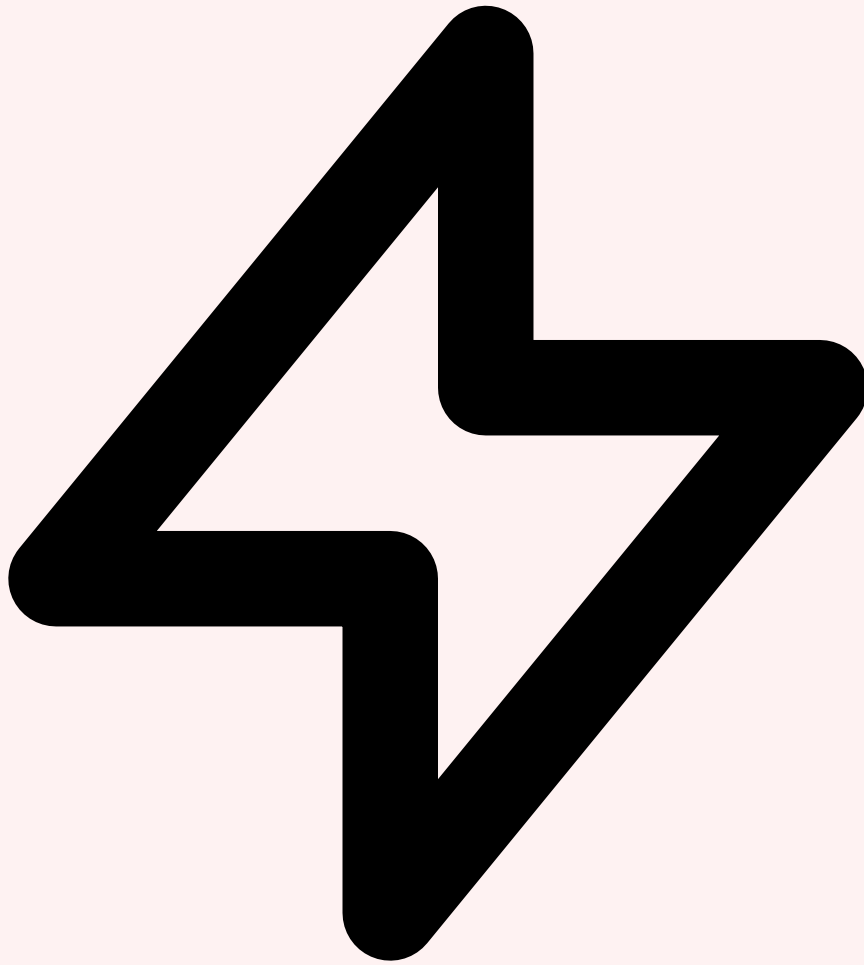
Llama 4, Mistral Large, Gemma 3 : Comparatif LLM Open Source constitue un enjeu majeur pour les professionnels de la sécurité informatique et les équipes techniques. Ce guide détaillé sur le comparatif LLM open source propose une méthodologie structurée, des outils éprouvés et des recommandations opérationnelles directement applicables. L'objectif est de fournir aux praticiens — consultants, ingénieurs sécurité, administrateurs systèmes — les connaissances et les techniques nécessaires pour aborder ce sujet avec rigueur. Chaque section s'appuie sur des retours d'expérience terrain et intègre les évolutions les plus récentes du domaine. Les recommandations présentées sont adaptées aux environnements d'entreprise et tiennent compte des contraintes opérationnelles réelles.

Table des Matières

1. [Le Paysage LLM Open Source en 2026](#)
2. [Llama 4 : Scout et Maverick par Meta](#)
3. [Mistral Large 2, Codestral et Pixtral](#)
4. [Gemma 3 : La Puissance Google en Open Source](#)
5. [Qwen 2.5 et DeepSeek V3 : Les Modèles Chinois](#)
6. [Benchmarks Comparatifs et Tableau Récapitulatif](#)
7. [Guide de Choix par Cas d'Usage](#)

Notre avis d'expert

La dynamique concurrentielle s'est considérablement intensifiée. **Meta, Google, Mistral AI, Alibaba et DeepSeek** se livrent une course à l'innovation qui profite directement aux entreprises et aux développeurs. Chaque trimestre apporte son lot de percées architecturales, de nouveaux records sur les benchmarks et d'optimisations qui repoussent les limites du possible sur du matériel accessible. Comparatif détaillé des LLM open source 2026 : Llama 4, Mistral Large, Gemma 3, Qwen 2.5, DeepSeek V3. Benchmarks, coûts et guide de choix. Guide. Dans un contexte où l'intelligence artificielle transforme les pratiques de cybersécurité, la maîtrise de ia comparatif llm open source devient un avantage stratégique pour les équipes techniques. Nous abordons notamment : table des matières, 2 llama 4 : scout et maverick par meta et 3 mistral large 2, codestral et pixtral. Les professionnels y trouveront des recommandations actionnables, des commandes prêtes à l'emploi et des stratégies de mise en œuvre adaptées aux environnements d'entreprise.

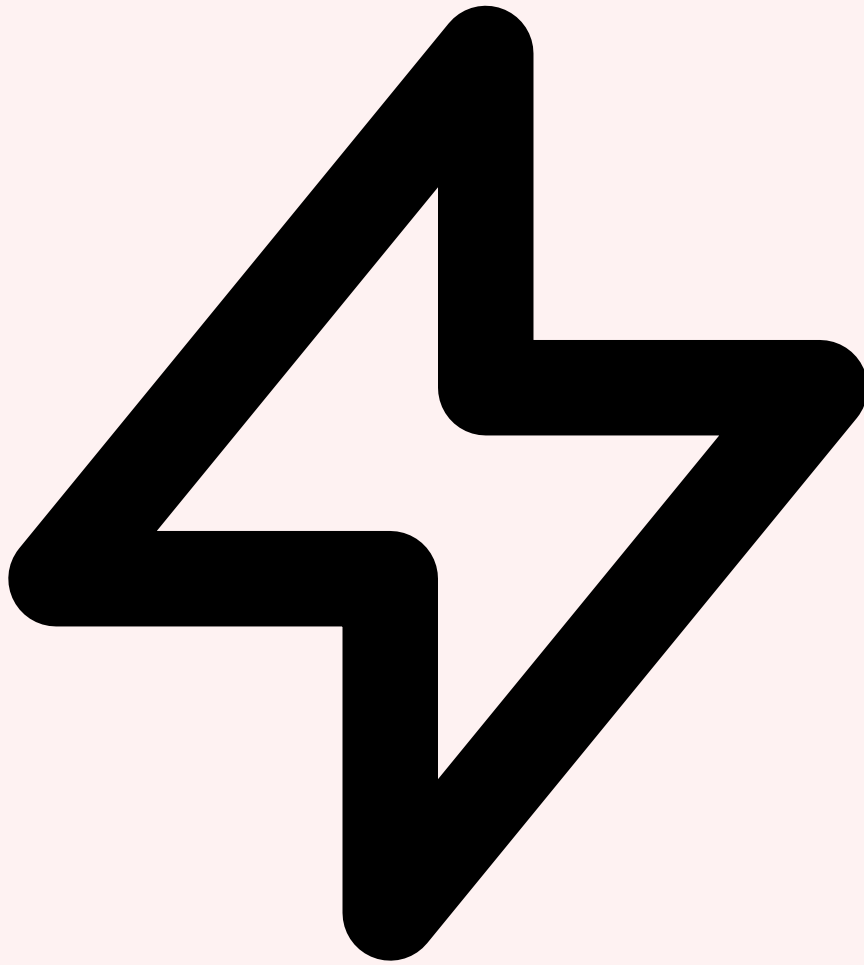


Une révolution en trois axes

Trois tendances structurantes redéfinissent le paysage LLM open source en 2026. Premièrement, l'architecture **Mixture of Experts (MoE)** s'est imposée comme le standard pour les modèles de grande taille. Llama 4, Mistral Large et DeepSeek V3 l'adoptent tous, permettant de multiplier la capacité du modèle sans multiplier proportionnellement le coût d'inférence. Un modèle de 400 milliards de paramètres totaux n'active typiquement que 50 à 100 milliards de paramètres par requête.

Deuxièmement, la **multimodalité native** est devenue la norme plutôt que l'exception. Les modèles de 2026 comprennent nativement texte, images, code et données structurées, ouvrant la porte à des applications qui étaient auparavant réservées aux API propriétaires comme GPT-4o ou Claude.

Troisièmement, les **fenêtres de contexte** ont explosé. Là où 4 096 tokens étaient la norme en 2023, on parle désormais de 128K à 10 millions de tokens, transformant radicalement les possibilités d'analyse documentaire, de raisonnement sur de longs textes et de génération augmentée par récupération (RAG).



Chronologie des releases majeures

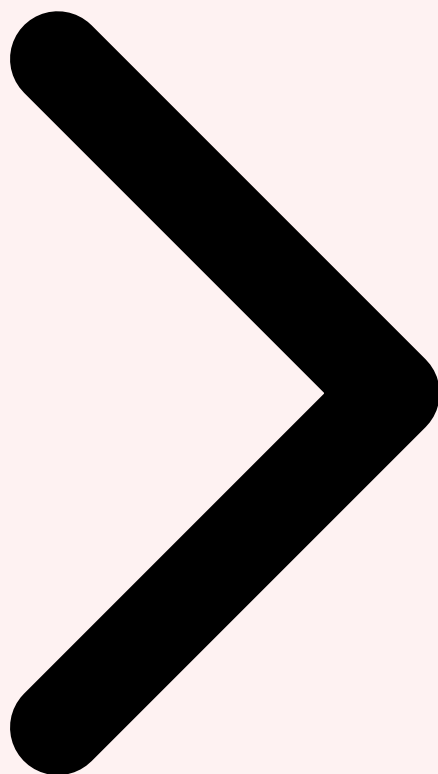
Pour bien comprendre l'accélération du rythme d'innovation, voici la chronologie des sorties majeures depuis fin 2024 jusqu'à début 2026. Chaque release a marqué une étape significative dans la démocratisation des LLM performants.

Cette chronologie illustre l'accélération impressionnante du rythme de publication. En l'espace de dix-huit mois, chaque acteur majeur a publié au moins une mise à jour significative, créant un écosystème en perpétuelle évolution. Pour les entreprises qui souhaitent adopter un LLM open source, le choix est à la fois plus riche et plus complexe que jamais.

Ce comparatif exhaustif passe en revue les cinq familles de modèles les plus pertinentes pour un déploiement professionnel en 2026. Pour chacune, nous analysons l'architecture technique, les performances sur les benchmarks standard, les cas d'usage privilégiés et les contraintes matérielles à anticiper. L'objectif est de vous fournir toutes les clés pour faire un choix éclairé, adapté à vos besoins spécifiques.



Table des Matières Paysage LLM 2026 Llama 4 (Meta)

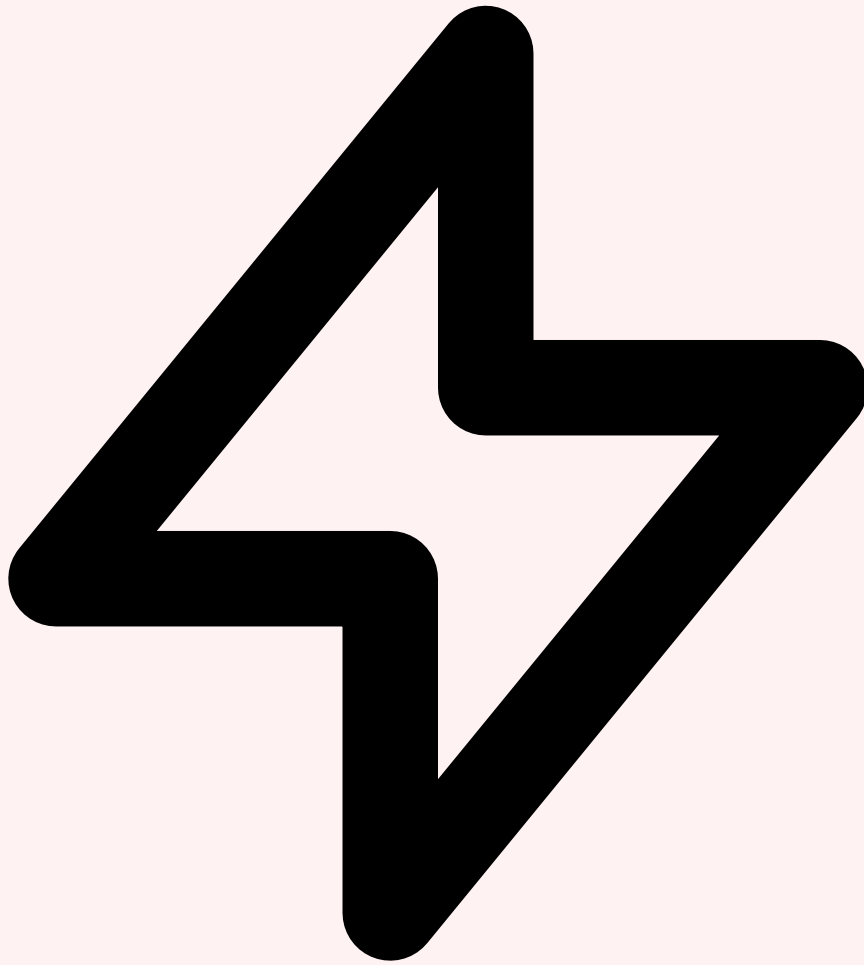


Cas concret

L'attaque par prompt injection sur les systèmes GPT documentée par OWASP en 2023 a révélé que des instructions malveillantes dissimulées dans des documents pouvaient détourner le comportement de chatbots d'entreprise, accédant à des données internes sensibles sans aucune authentification supplémentaire.

2 Llama 4 : Scout et Maverick par Meta

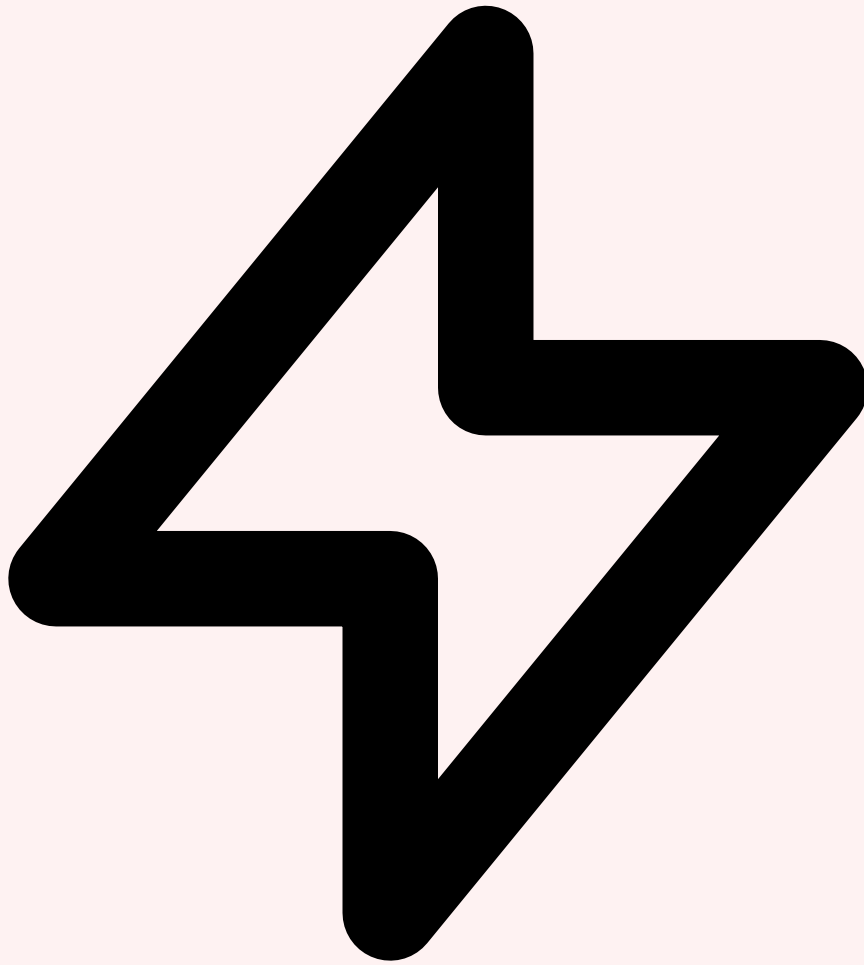
Avec **Llama 4**, Meta franchit un cap majeur dans sa stratégie open source. La quatrième génération de sa famille de modèles introduit pour la première fois une **architecture Mixture of Experts (MoE)** qui change fondamentalement l'équation performance-efficacité. Deux variantes principales sont disponibles : **Scout**, optimisé pour le déploiement sur une seule machine, et **Maverick**, la version premium taillée pour les charges de travail les plus exigeantes.



Architecture MoE : le saut technologique

Llama 4 Scout embarque 109 milliards de paramètres au total, organisés en 16 experts dont seulement 2 sont activés par token. Cela signifie que le coût d'inférence effectif correspond à environ 17 milliards de paramètres actifs — une efficacité remarquable qui permet au modèle de tourner sur un unique serveur équipé d'un GPU H100 80GB. La fenêtre de contexte native atteint **10 millions de tokens**, un record absolu pour un modèle de cette taille, ouvrant des possibilités inédites en analyse documentaire massive. Pour approfondir, consultez [IA Offensive : Comment les Attaquants Utilisent les LLM](#).

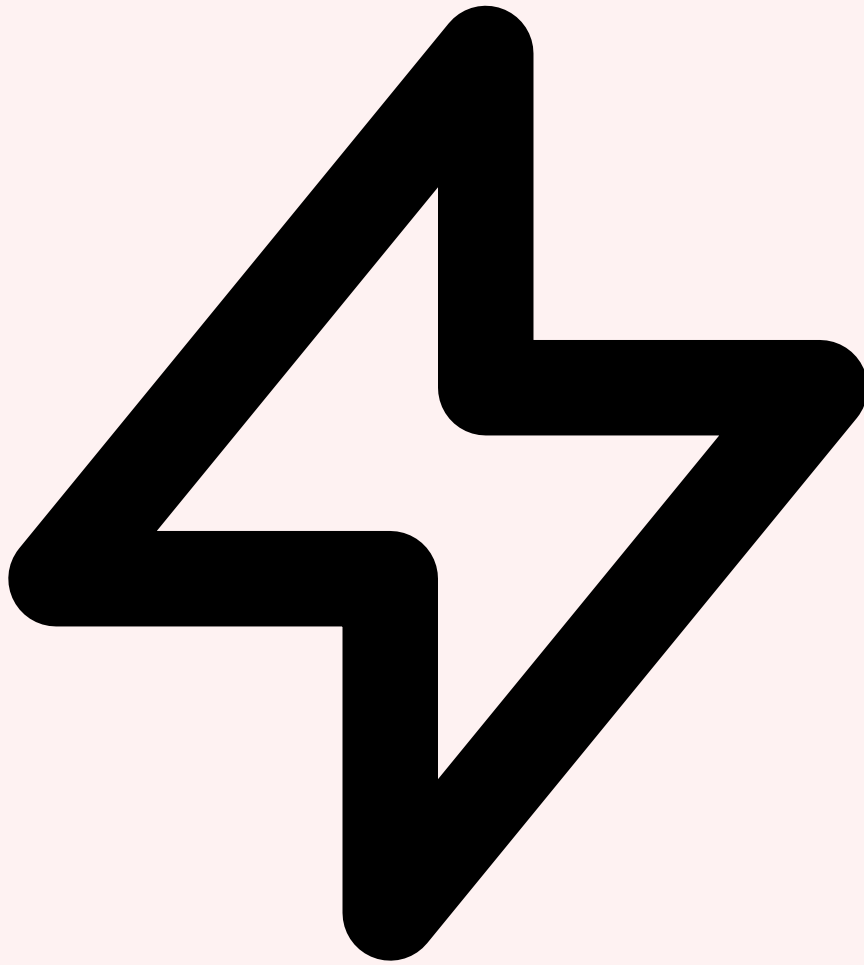
Llama 4 Maverick monte en puissance avec 400 milliards de paramètres totaux, 128 experts et une activation de 17 milliards de paramètres par token. Ce modèle cible les entreprises ayant besoin de la meilleure qualité possible sur des tâches complexes : raisonnement multi-étapes, génération de code complexe, analyse juridique ou médicale. Sa fenêtre de contexte de **1 million de tokens** reste exceptionnelle pour un modèle de cette envergure.



Performances et benchmarks

Sur les benchmarks standard, Llama 4 affiche des résultats qui le placent systématiquement dans le top 3 des modèles open source. Scout obtient un score **MMLU de 85.4%**, surpassant Llama 3.1 70B de plus de 3 points. Sur HumanEval (génération de code), il atteint **84.2%**, démontrant une maîtrise solide du code dans plus de 20 langages de programmation. Maverick pousse encore plus loin avec **89.3% sur MMLU** et **88.1% sur HumanEval**, rivalisant directement avec GPT-4o sur ces métriques.

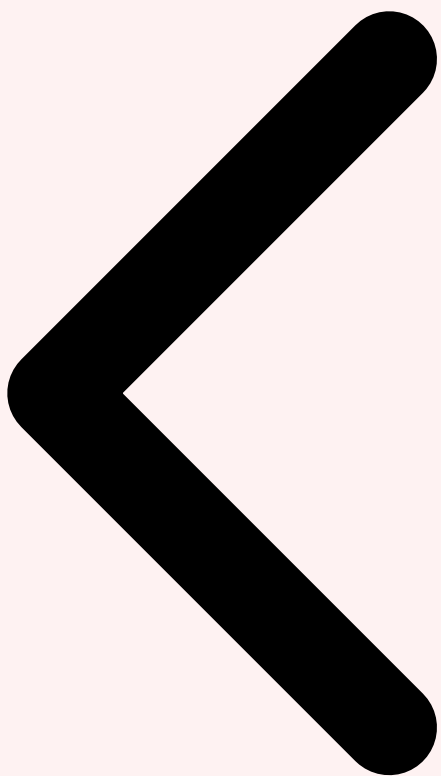
La **multimodalité** est intégrée nativement dans les deux variantes. Llama 4 comprend les images avec des performances de pointe sur les benchmarks visuels (MMMU, ChartQA, DocVQA), ce qui en fait un outil polyvalent pour l'analyse de documents, la compréhension de schémas techniques ou l'extraction d'informations à partir de captures d'écran.



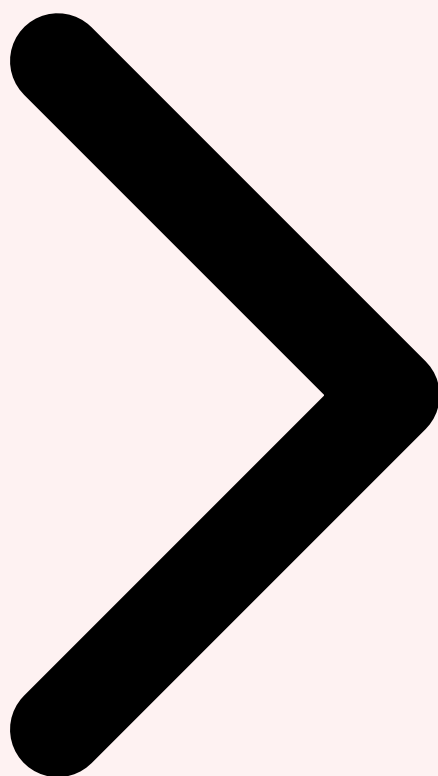
Déploiement et licence

Llama 4 est distribué sous la **Llama Community License**, qui autorise l'utilisation commerciale pour les organisations de moins de 700 millions d'utilisateurs actifs mensuels. Les poids sont disponibles sur Hugging Face et le déploiement est supporté nativement par vLLM, TGI, Ollama et llama.cpp. Pour Scout en quantization INT4, comptez environ **32 Go de VRAM** — accessible sur un RTX 4090 ou un A6000.

- **Points forts** — Architecture MoE efficace, contexte jusqu'à 10M tokens, multimodalité native, excellente performance par dollar
- **Limites** — Licence restrictive au-delà de 700M utilisateurs, Maverick nécessite un cluster multi-GPU, fine-tuning MoE plus complexe
- **Cas d'usage idéaux** — Chatbots d'entreprise, analyse documentaire, génération de code, RAG sur corpus volumineux



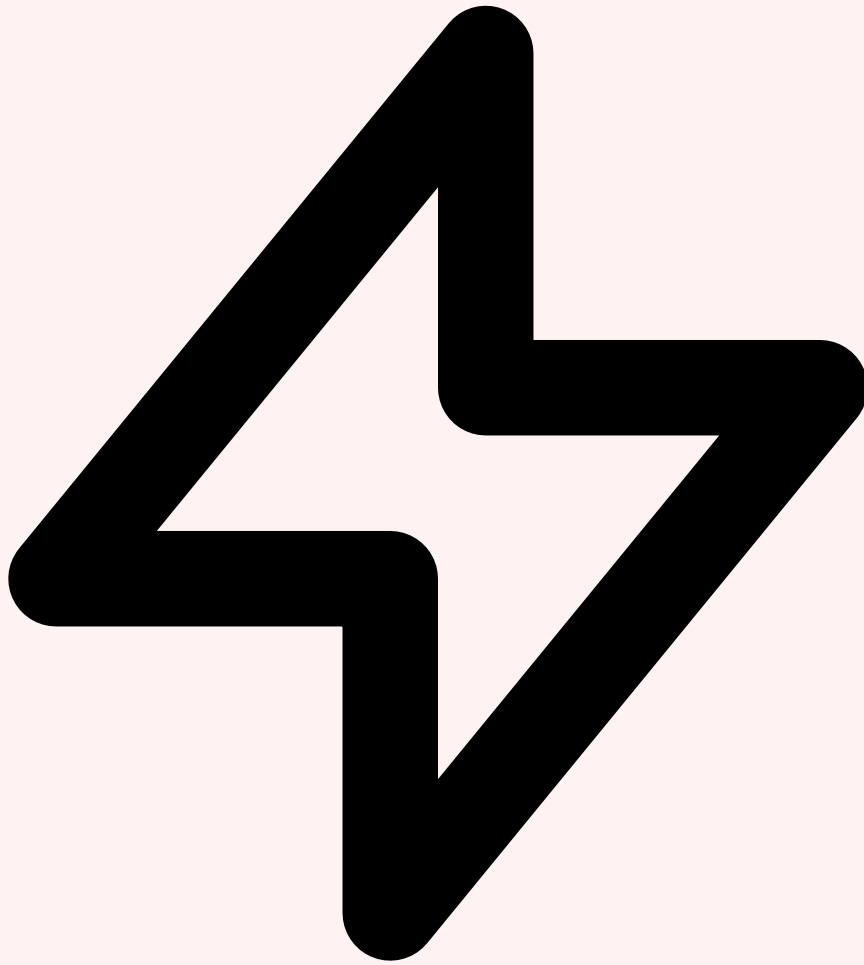
Paysage LLM 2026 Llama 4 (Meta) Mistral Large 2



Vos pipelines de données d'entraînement sont-ils protégés contre l'empoisonnement ?

3 Mistral Large 2, Codestral et Pixtral

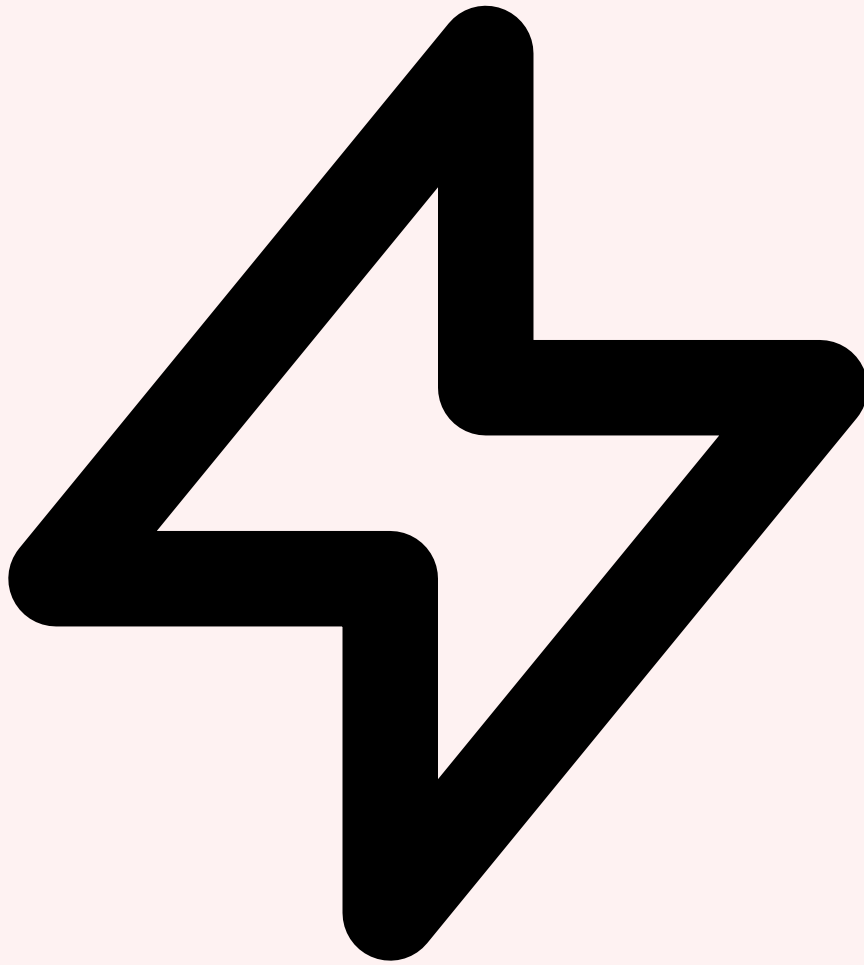
Mistral AI, la pépite française fondée par d'anciens chercheurs de Meta et Google DeepMind, s'est imposée comme un acteur incontournable de l'écosystème LLM européen. Avec **Mistral Large 2**, l'entreprise propose un modèle dense de 123 milliards de paramètres qui se distingue par sa maîtrise exceptionnelle du français et des langues européennes, un atout différenciant majeur pour les entreprises francophones.



Mistral Large 2 : le modèle généraliste

Mistral Large 2 est un modèle **dense de 123B paramètres** avec une fenêtre de contexte de **128K tokens**. Contrairement aux approches MoE de Llama 4 ou DeepSeek, Mistral opte pour une architecture dense optimisée, arguant que la stabilité de l'entraînement et la prédictibilité des performances justifient le surcoût en inférence. Le modèle supporte nativement le **function calling** structuré et le mode JSON, facilitant son intégration dans des pipelines d'agents IA.

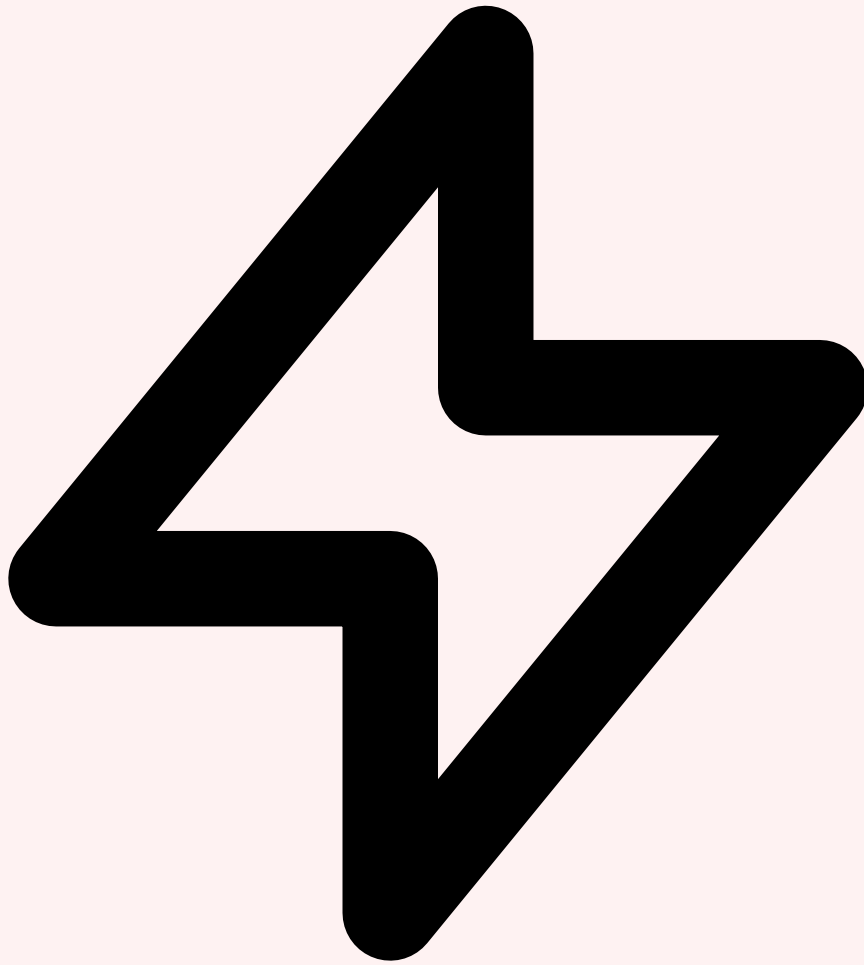
Sur les benchmarks multilingues, Mistral Large 2 excelle particulièrement. Il obtient un score **MMLU de 84.0%** en moyenne, mais atteint **87.2% sur les sous-ensembles en français**, surpassant tous les concurrents sur ce critère. Son entraînement a bénéficié de données de haute qualité en français, allemand, espagnol et italien, un avantage stratégique pour les déploiements européens soumis aux contraintes du RGPD.



Codestral : le spécialiste du code

Codestral est le modèle dédié à la génération et à la compréhension de code de Mistral AI. Basé sur une architecture de 22 milliards de paramètres optimisée pour la latence, il supporte plus de **80 langages de programmation** et se distingue par sa capacité à générer du code idiomatique et bien structuré. Sur HumanEval, Codestral atteint **86.5%**, le plaçant au niveau des meilleurs modèles spécialisés comme Code Llama et DeepSeek Coder V2.

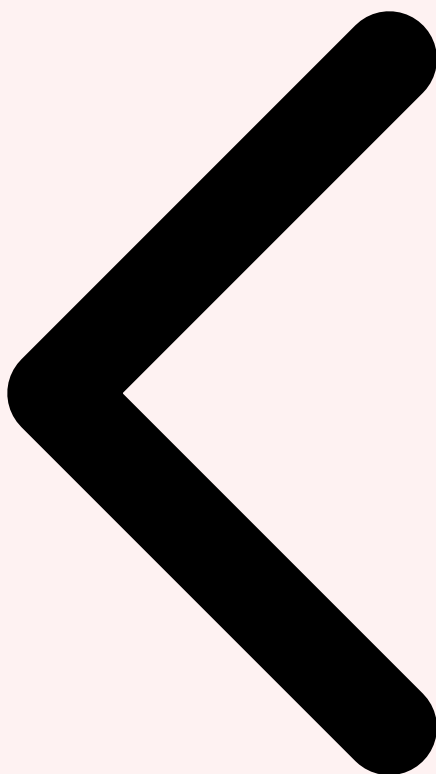
L'intégration native dans les IDE (VS Code, JetBrains, Neovim) via le protocole Continue et la compatibilité avec le format OpenAI en font un excellent candidat pour remplacer GitHub Copilot dans les environnements où la souveraineté des données est critique. Le modèle est disponible sous **licence non-commerciale** pour la recherche, et sous licence commerciale via l'API Mistral.



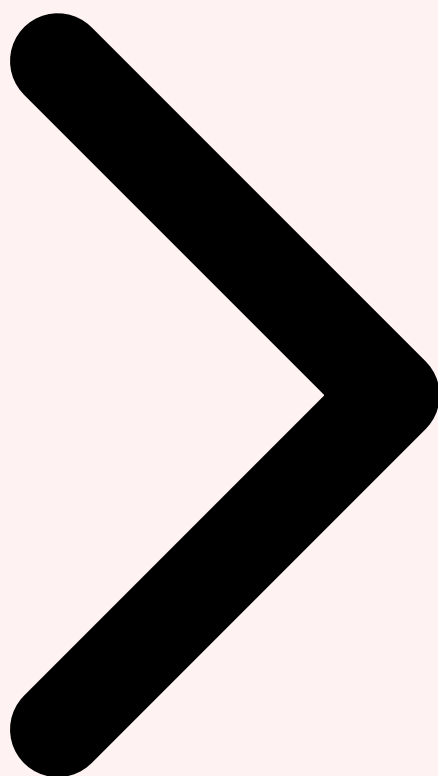
Pixtral : la vision multimodale

Pixtral Large complète l'offre Mistral avec un modèle multimodal de 124 milliards de paramètres capable de comprendre texte et images simultanément. Son architecture combine un encodeur vision de 400M de paramètres avec le backbone Mistral Large, permettant l'analyse de documents complexes, de graphiques et de captures d'écran. Pixtral atteint des scores de pointe sur **DocVQA (93.2%)** et **ChartQA (88.4%)**, rivaux des meilleurs modèles propriétaires.

- **Points forts** — Excellence en français et langues européennes, fonction calling natif, écosystème complet (code + vision), conformité RGPD
- **Limites** — Architecture dense plus coûteuse en inférence, Codestral sous licence restrictive, contexte limité à 128K vs 10M pour Llama 4
- **Cas d'usage idéaux** — Entreprises francophones, conformité RGPD, assistant code souverain, analyse documentaire européenne

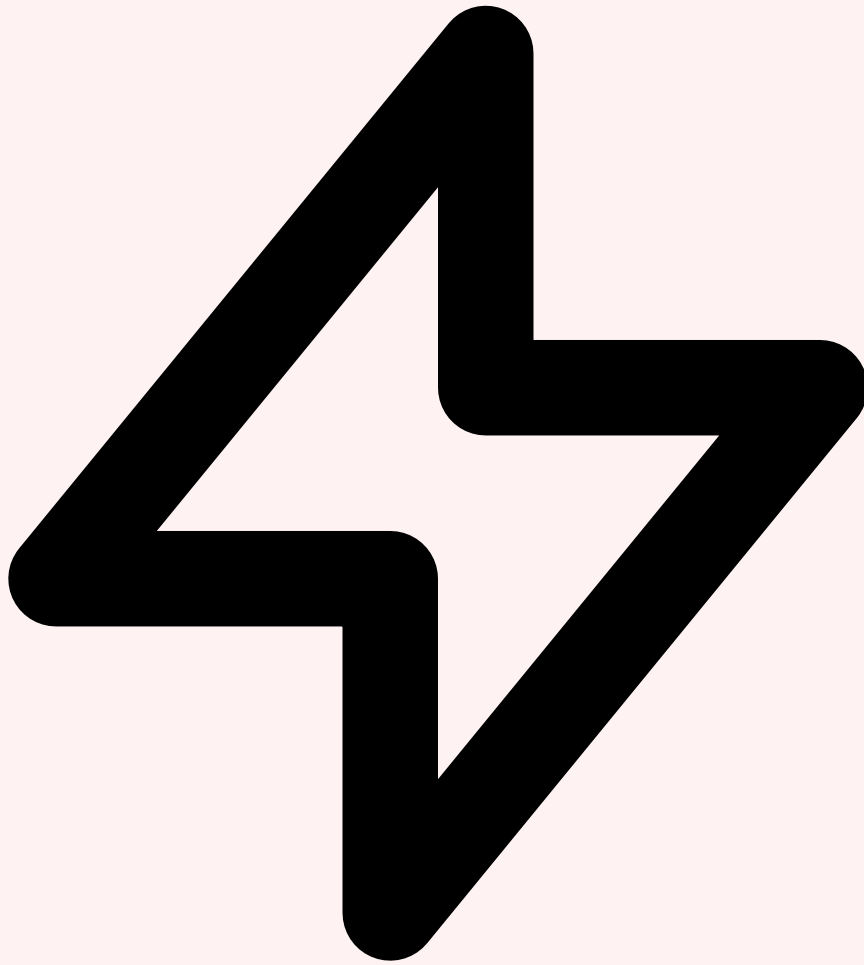


Llama 4 (Meta) Mistral Large 2 Gemma 3 (Google)



4 Gemma 3 : La Puissance Google en Open Source

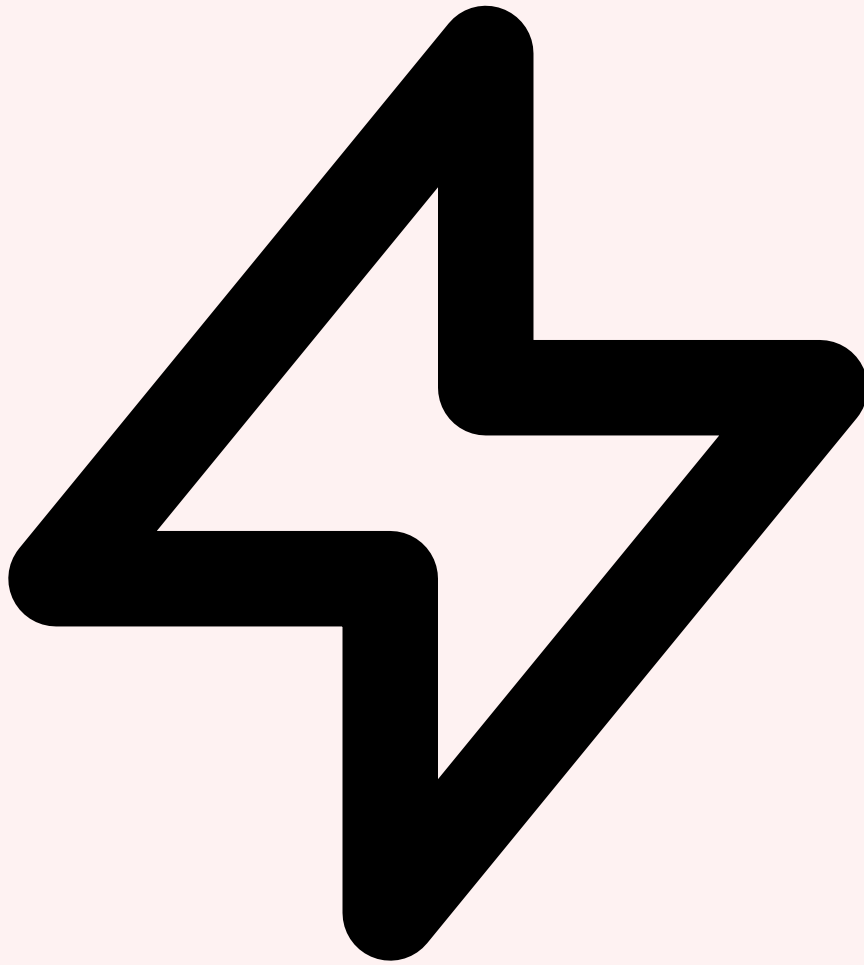
Gemma 3 représente la troisième génération de la famille de modèles open source de Google DeepMind. Construite sur les fondations de l'architecture Gemini, cette famille se distingue par une gamme de tailles exceptionnellement étendue — de **1B à 27B paramètres** — qui couvre l'intégralité du spectre, du déploiement sur smartphone jusqu'au serveur d'entreprise. C'est cette polyvalence qui fait de Gemma 3 un choix stratégique pour les organisations qui ont besoin de déployer le même modèle à différentes échelles. Pour approfondir, consultez [Small Language Models : Phi-4, Gemma et IA Embarquée](#).



Architecture et déclinaisons

Gemma 3 se décline en quatre tailles principales : **1B, 4B, 12B et 27B** paramètres. Chaque taille est disponible en version pré-entraînée (PT) et en version instruction-tuned (IT). Le modèle 27B est le porte-étendard de la famille, offrant des performances qui rivalisent avec des modèles deux à trois fois plus grands grâce aux optimisations héritées de l'entraînement de Gemini.

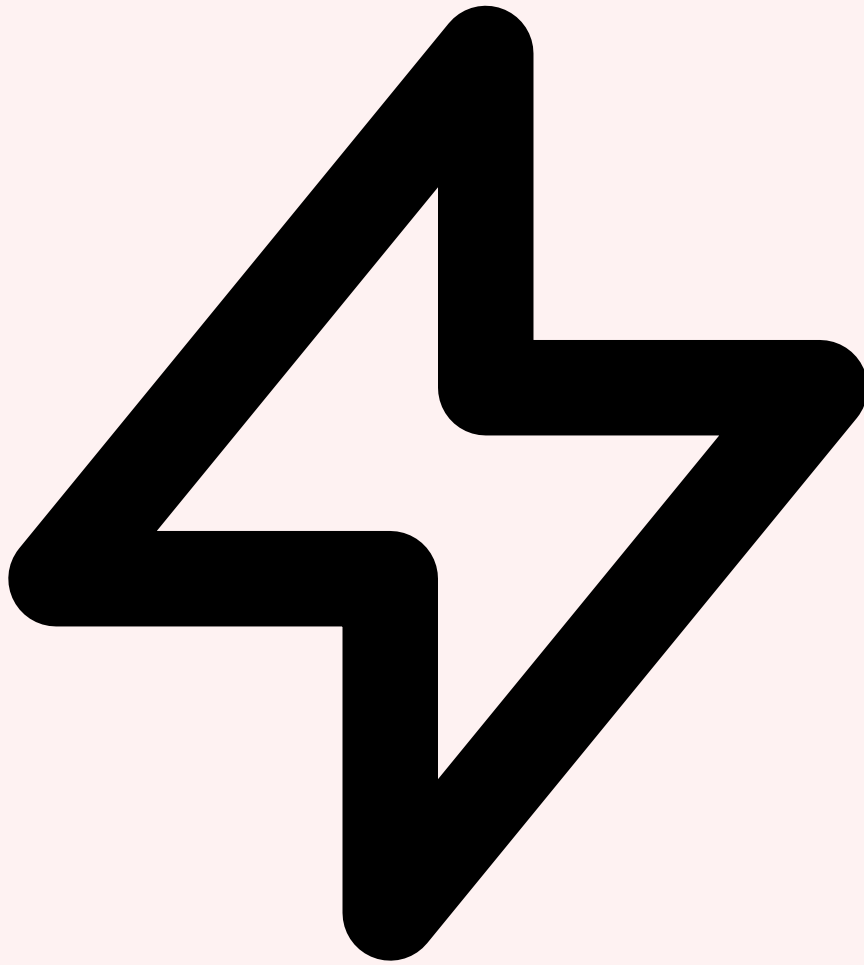
L'architecture de Gemma 3 intègre plusieurs innovations notables. Le **sliding window attention** alterne avec l'attention globale pour optimiser l'utilisation mémoire sur les longues séquences. Le modèle 27B supporte une fenêtre de contexte de **128K tokens**, comparable à Mistral Large 2. La multimodalité est native à partir de la taille 4B, avec un encodeur vision SigLIP2 capable de traiter des images haute résolution.



ShieldGemma : la sécurité intégrée

Un différenciateur majeur de l'écosystème Gemma est **ShieldGemma**, un ensemble de modèles de garde (guardrails) spécialisés dans la détection de contenus dangereux, toxiques ou inappropriés. ShieldGemma 2 fonctionne comme un filtre de sécurité multimodal capable d'analyser à la fois le texte et les images pour détecter les violations de politique de contenu. Cette approche de la sécurité IA intégrée est particulièrement valorisée dans les secteurs réglementés comme la santé, la finance et l'éducation.

Google a également publié **Gemma Scope**, un outil d'interprétabilité qui utilise des autoencodeurs sparse pour comprendre les mécanismes internes du modèle. Cette transparence est un atout considérable pour les organisations qui doivent justifier les décisions prises par leur IA auprès de régulateurs ou d'auditeurs.

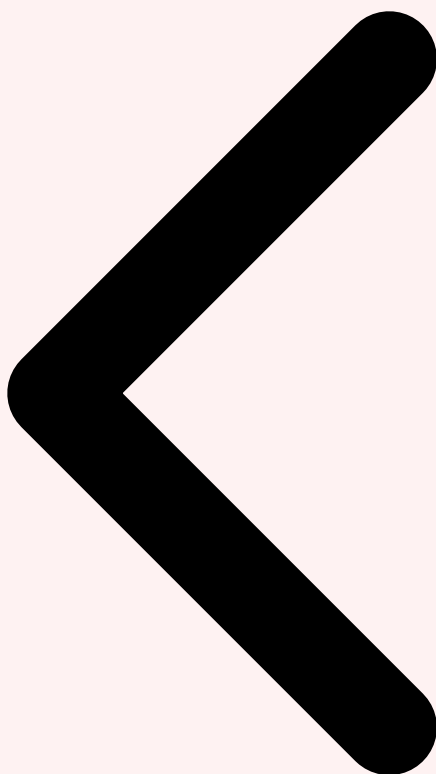


Optimisation mobile et edge

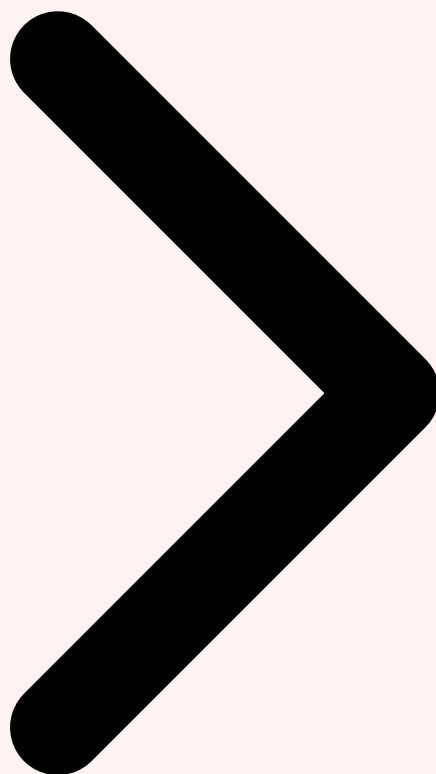
C'est sur le segment **mobile et edge computing** que Gemma 3 brille particulièrement. Le modèle 4B quantifié en INT4 fonctionne confortablement sur un smartphone Android haut de gamme avec seulement **3 Go de RAM**. Les optimisations spécifiques pour les processeurs ARM et les GPU mobiles (Adreno, Mali) garantissent des temps de réponse acceptables même sans connexion réseau.

Le modèle 1B, quant à lui, peut tourner sur des appareils IoT et des systèmes embarqués avec des contraintes mémoire extrêmes. Google a démontré son déploiement sur des Raspberry Pi 5 et des Jetson Nano, ouvrant la voie à des applications IA véritablement décentralisées dans l'industrie, l'agriculture intelligente ou la domotique.

- **Points forts** — Gamme de tailles complète (1B-27B), ShieldGemma pour la sécurité, optimisation mobile/edge de référence, licence permissive
- **Limites** — Taille maximale de 27B (pas de modèle 70B+), performances brutes inférieures à Llama 4 Maverick, communauté de fine-tuning moins développée
- **Cas d'usage idéaux** — Déploiement mobile/edge, applications embarquées, IA sécurisée dans les secteurs réglementés, prototypage rapide

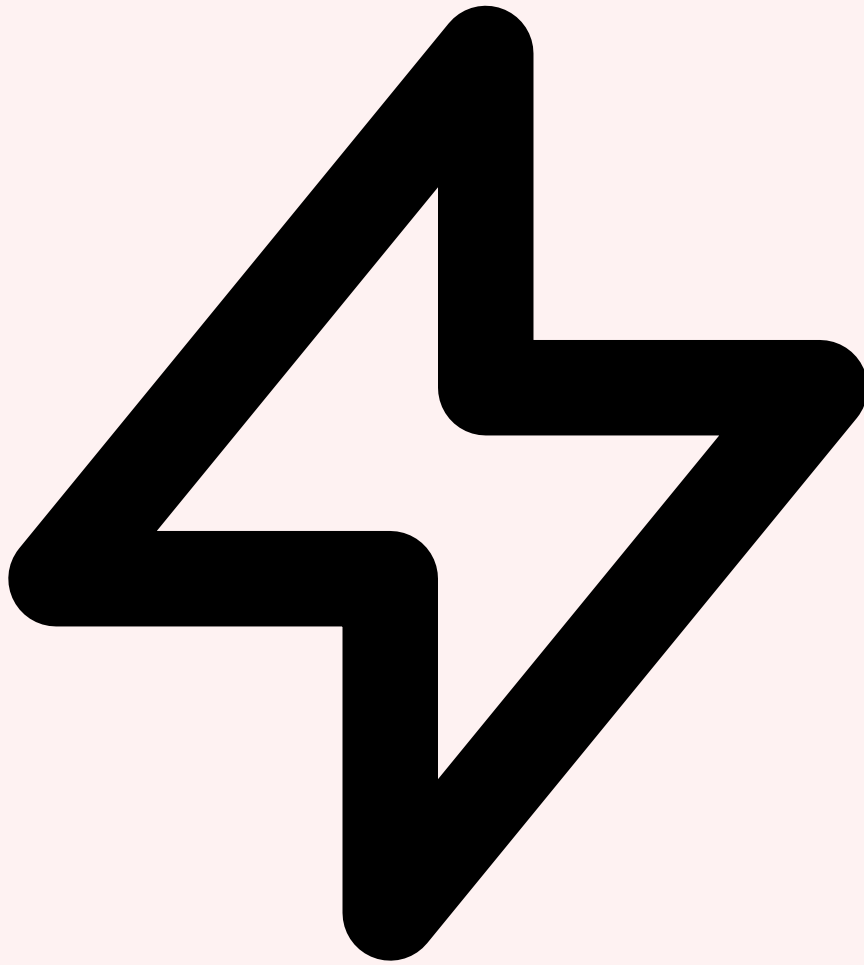


Mistral Large 2 Gemma 3 (Google) Qwen & DeepSeek



5 Qwen 2.5 et DeepSeek V3 : Les Modèles Chinois

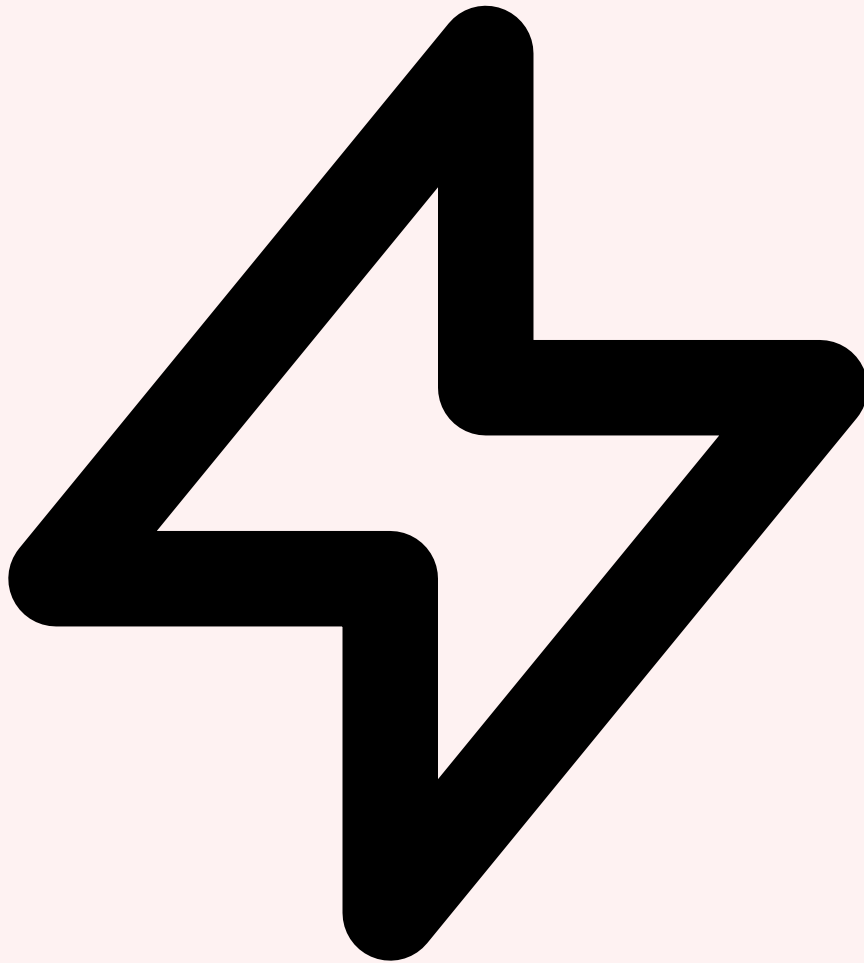
L'émergence des modèles open source chinois a été l'une des surprises majeures de la période 2024-2026. **Qwen 2.5 d'Alibaba** et **DeepSeek V3** ont démontré que l'innovation en matière de LLM n'est plus l'apanage exclusif des laboratoires américains. Ces modèles rivalisent frontalement avec les meilleurs modèles occidentaux sur les benchmarks internationaux, tout en offrant des performances exceptionnelles sur le chinois et les langues asiatiques.



Qwen 2.5 : la gamme complète d'Alibaba

Qwen 2.5 se décline en une gamme impressionnante : 0.5B, 1.5B, 3B, 7B, 14B, 32B et 72B paramètres. Le modèle phare de **72 milliards de paramètres** est celui qui retient le plus l'attention pour les déploiements professionnels. Entraîné sur 18 000 milliards de tokens couvrant 29 langues, il affiche un score MMLU de **85.3%** et un HumanEval de **86.4%**, des performances qui le placent au niveau de Llama 4 Scout.

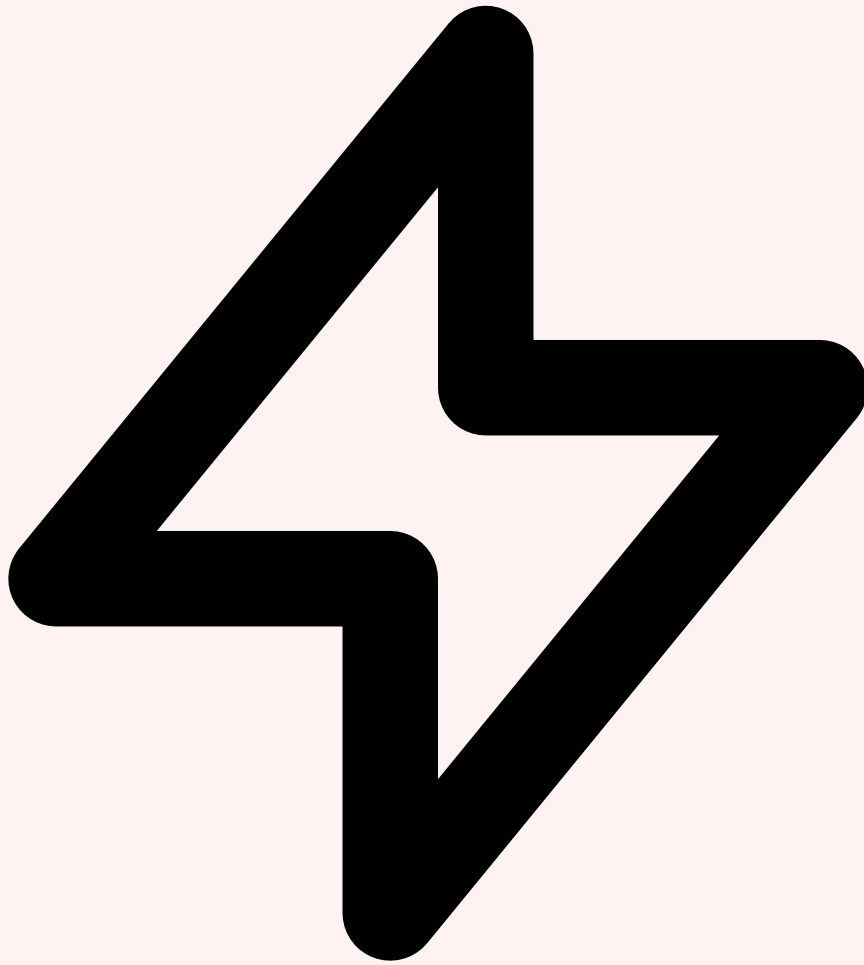
Alibaba a également publié des variantes spécialisées qui enrichissent l'écosystème. **Qwen 2.5-Coder** est optimisé pour la génération et la compréhension de code, avec un support de 92 langages de programmation. **Qwen 2.5-Math** excelle en raisonnement mathématique, surpassant GPT-4o sur les benchmarks MATH et GSM8K. Enfin, **Qwen-VL** offre des capacités multimodales compétitives pour l'analyse d'images et de vidéos.



DeepSeek V3 : l'efficacité radicale

DeepSeek V3 a fait sensation en janvier 2025 avec une approche qui a redéfini les standards d'efficacité. Ce modèle MoE de **671 milliards de paramètres totaux** (37B actifs par token) a été entraîné pour un coût estimé de seulement **5.6 millions de dollars** — une fraction du budget des modèles comparables. Cette prouesse repose sur des innovations architecturales comme le Multi-Head Latent Attention (MLA) et le DeepSeekMoE avec routage auxiliaire-free.

Les performances de DeepSeek V3 sont remarquables : **87.1% sur MMLU**, **89.2% sur HumanEval** et des résultats de pointe sur les benchmarks mathématiques. Le modèle excelle particulièrement en **raisonnement et en code**, domaines où il rivalise avec Claude 3.5 Sonnet et GPT-4o. Son successeur, DeepSeek-R1, a introduit le approche du raisonnement par chaîne de pensée (chain-of-thought) avec des résultats exceptionnels sur les problèmes complexes. Pour approfondir, consultez [Fine-Tuning de LLM Open Source : Guide Complet LoRA et QLoRA](#).



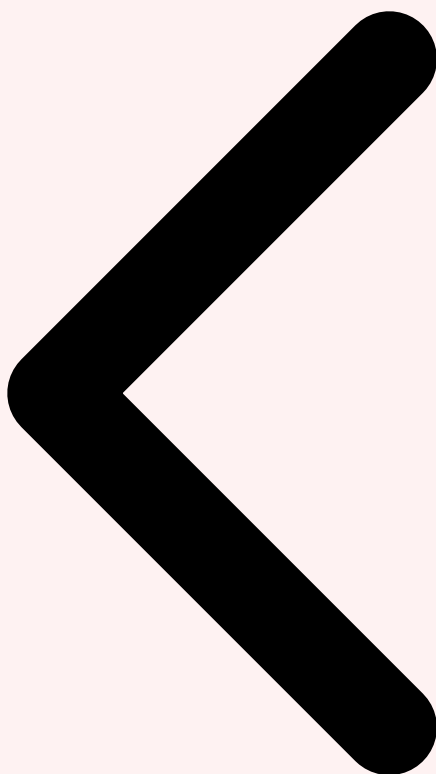
Considérations géopolitiques et pratiques

L'adoption des modèles chinois en Europe soulève des questions légitimes. Sur le plan **technique**, les deux modèles sont distribués sous des licences permissives (Apache 2.0 pour Qwen, MIT-like pour DeepSeek) et les poids sont intégralement disponibles sur Hugging Face. Cependant, certaines organisations expriment des réserves sur la **souveraineté des données d'entraînement** et les potentielles backdoors, même si aucune preuve concrète n'a été apportée à ce jour.

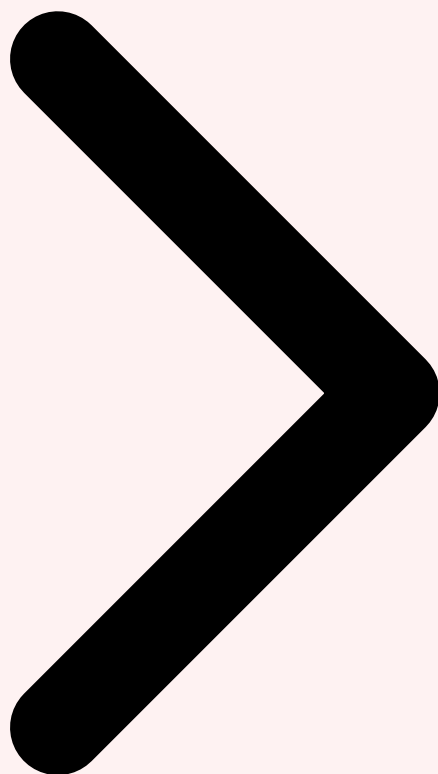
En pratique, le déploiement on-premise élimine les risques liés à l'exfiltration de données puisque le modèle tourne intégralement sur votre infrastructure. L'analyse des poids par la communauté open source n'a révélé aucun comportement suspect. Pour les organisations sensibles, une approche pragmatique consiste à **utiliser ces modèles pour les tâches non confidentielles** tout en réservant un modèle occidental audité pour les données sensibles.

- **Qwen 2.5 — Points forts** — Gamme complète, variantes spécialisées (code, math, vision), licence Apache 2.0, excellent rapport qualité/prix
- **DeepSeek V3 — Points forts** — Performances top-tier, coût d'entraînement transformateur, architecture MoE innovante, excellente en raisonnement

- **Limites communes** — Perception géopolitique, support communautaire principalement sinophone, documentation technique parfois incomplète en anglais

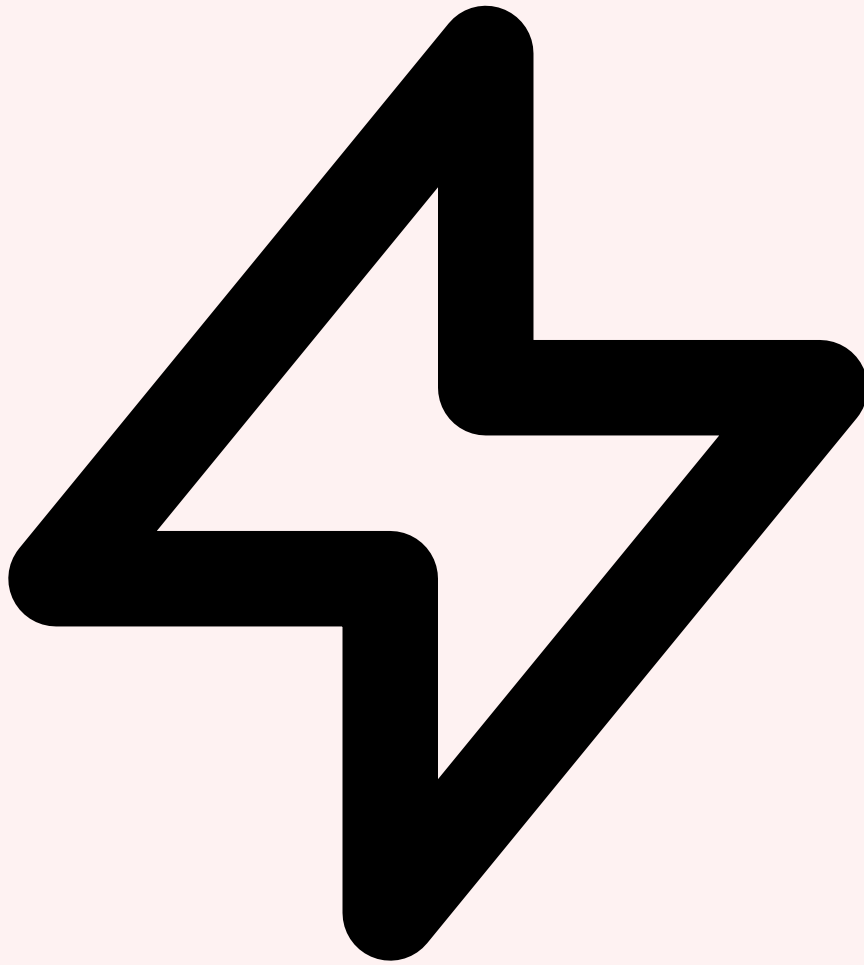


Gemma 3 (Google) Qwen & DeepSeek Benchmarks Comparatifs



6 Benchmarks Comparatifs et Tableau Récapitulatif

Comparer des LLM entre eux exige une méthodologie rigoureuse. Les benchmarks standardisés offrent un cadre objectif, mais ce que chaque métrique mesure réellement et quelles sont ses limites. Dans cette section, nous passons en revue les résultats consolidés des cinq familles de modèles sur les benchmarks les plus pertinents pour un déploiement professionnel.



Comprendre les métriques

Avant de plonger dans les chiffres, clarifions les principaux benchmarks utilisés dans ce comparatif :

- **▸MMLU (Massive Multitask Language Understanding)** — Évalue les connaissances générales et le raisonnement sur 57 domaines académiques (sciences, histoire, droit, médecine). Score en pourcentage, le seuil professionnel se situe au-dessus de 80%.
- **▸HumanEval** — Mesure la capacité de génération de code Python fonctionnel à partir de docstrings. 164 problèmes de programmation, score pass@1 en pourcentage.
- **▸MT-Bench** — Benchmark conversationnel multi-tour évalué par GPT-4. Score sur 10, mesure la qualité des réponses dans des conversations réalistes avec suivi de contexte.
- **▸MATH** — Problèmes mathématiques de niveau compétition. Évalue le raisonnement formel et la résolution de problèmes complexes étape par étape.

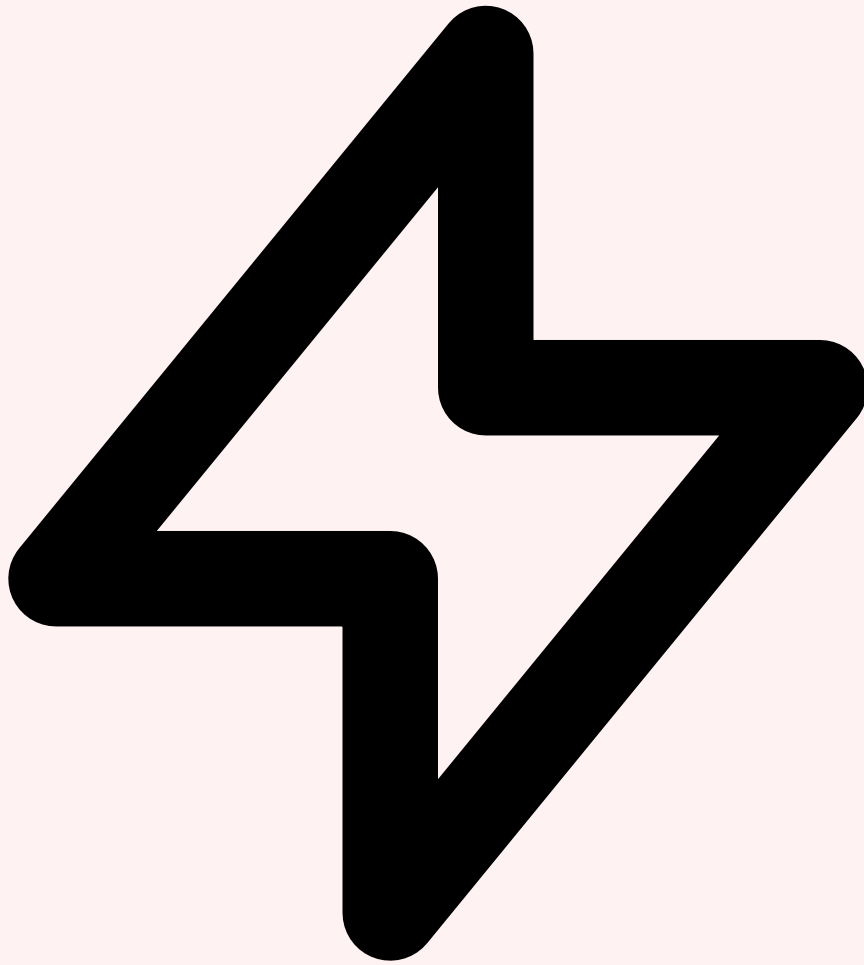
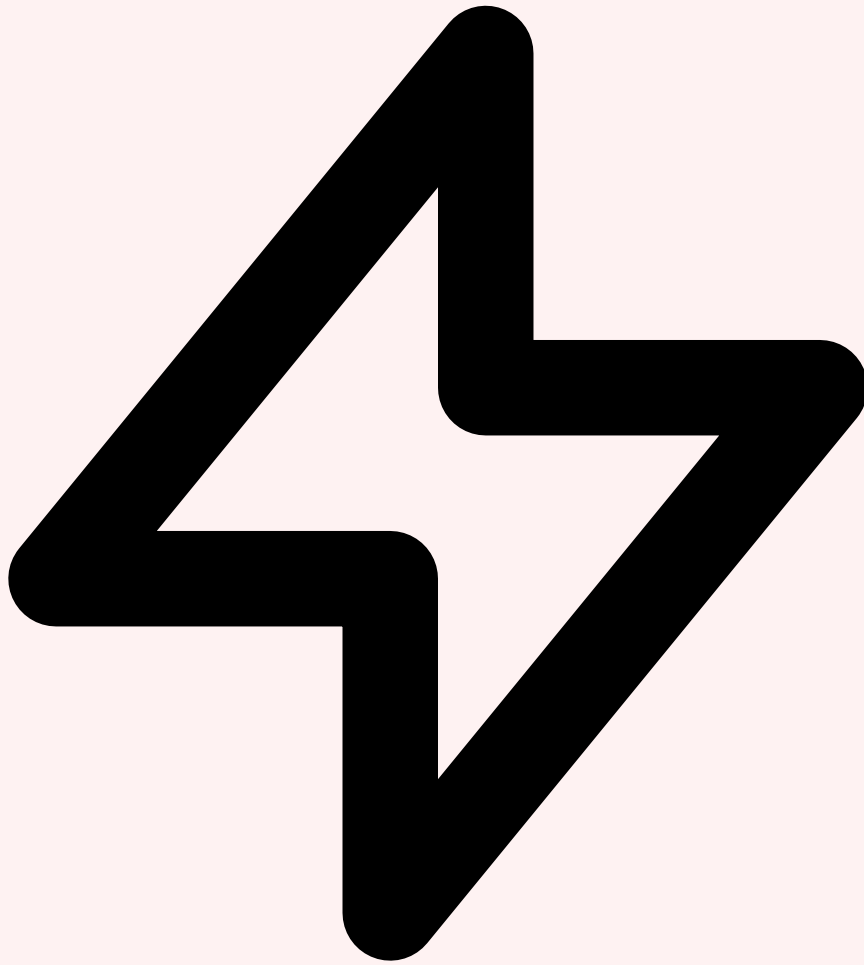


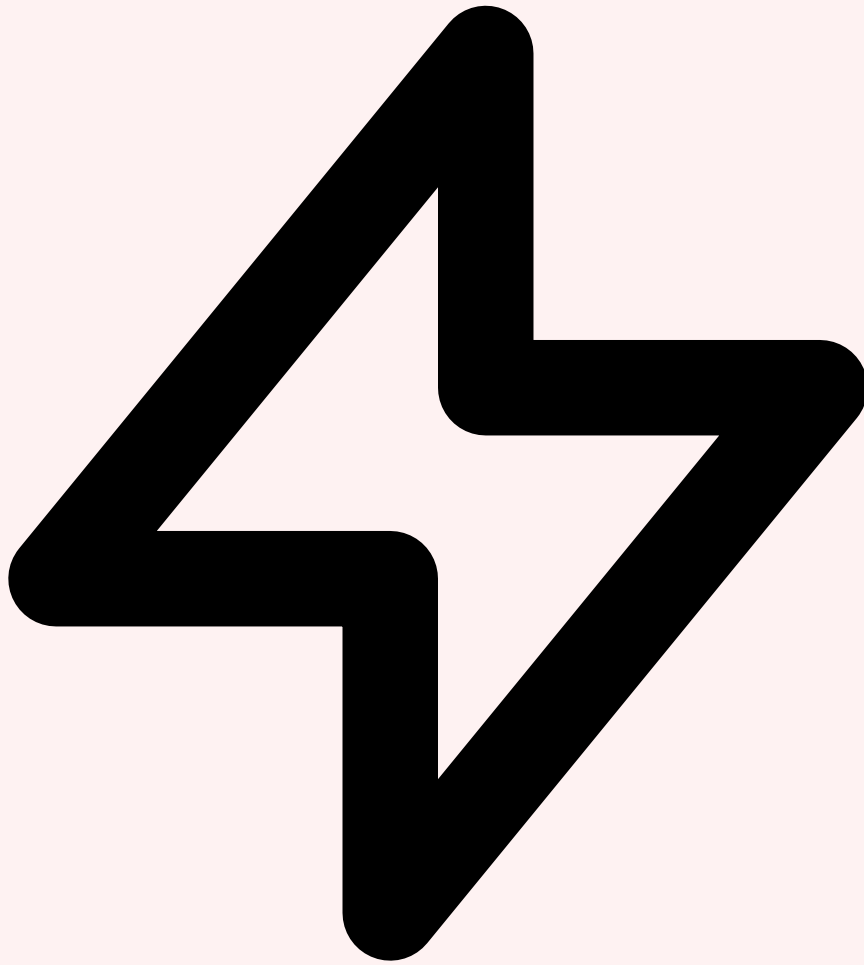
Tableau comparatif détaillé

Modèle	Params (actifs)	MMLU	HumanEval	MT-Bench	Contexte	Licence
Llama 4 Scout	109B (17B)	85.4%	84.2%	8.7	10M	Llama CL
Llama 4 Maverick	400B (17B)	89.3%	88.1%	9.1	1M	Llama CL
Mistral Large 2	123B (dense)	84.0%	81.9%	8.6	128K	Research
Gemma 3 27B	27B (dense)	78.7%	74.3%	8.3	128K	Permissive
Qwen 2.5 72B	72B (dense)	85.3%	86.4%	8.8	128K	Apache 2.0
DeepSeek V3	671B (37B)	87.1%	89.2%	9.0	128K	MIT-like



Analyse radar multi-dimensionnelle

Le tableau ci-dessus ne raconte qu'une partie de l'histoire. Pour une vision plus holistique, le diagramme radar ci-dessous compare les modèles sur six dimensions clés : connaissances générales (MMLU), code (HumanEval), conversation (MT-Bench), raisonnement mathématique (MATH), multimodalité et efficacité de déploiement.

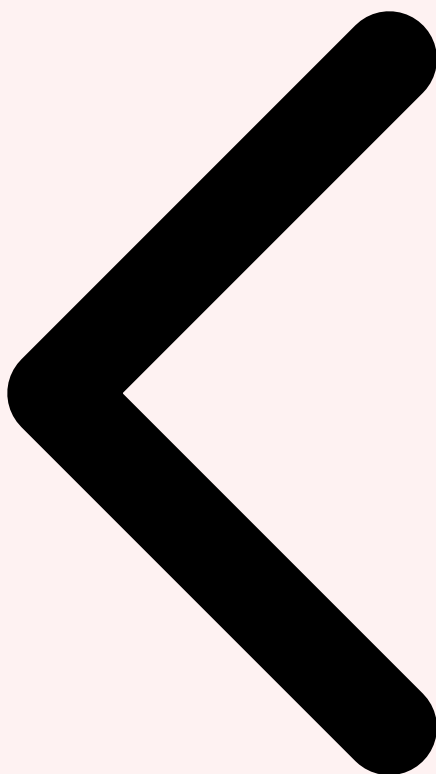


Interprétation des résultats

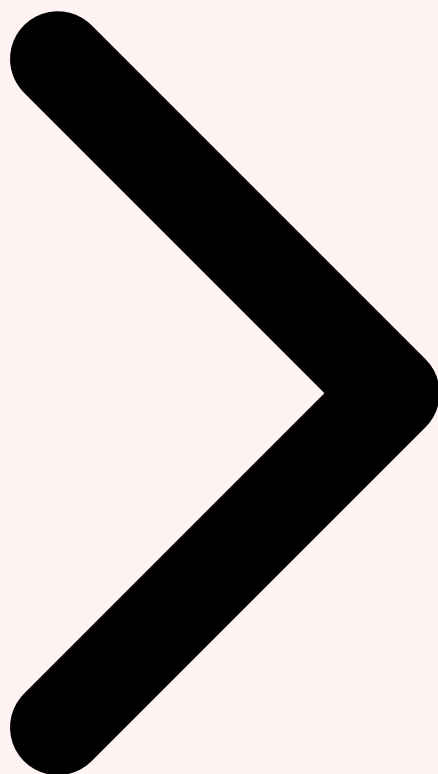
Plusieurs enseignements se dégagent de cette analyse comparative. **Llama 4 Maverick et DeepSeek V3** dominent le classement général, avec des performances quasi équivalentes sur la plupart des métriques. Le choix entre les deux dépendra principalement de vos contraintes de déploiement et de votre sensibilité géopolitique. Maverick a l'avantage de la fenêtre de contexte gigantesque (1M tokens), tandis que DeepSeek V3 impressionne par son efficacité architecturale.

Qwen 2.5 72B est la surprise de ce comparatif, offrant des performances proches des géants MoE avec un modèle dense plus simple à déployer et à fine-tuner. **Mistral Large 2** se démarque par son excellence linguistique en français et son écosystème européen. **Gemma 3 27B**, bien que moins performant en valeur absolue, offre le meilleur ratio performance/taille et reste imbattable sur le segment mobile et edge.

Il est crucial de rappeler que les benchmarks ne sont qu'un indicateur parmi d'autres. La performance réelle sur votre cas d'usage spécifique peut différer significativement des scores standardisés. Nous recommandons toujours de **tester les modèles candidats sur un échantillon représentatif de vos données réelles** avant de prendre une décision finale.

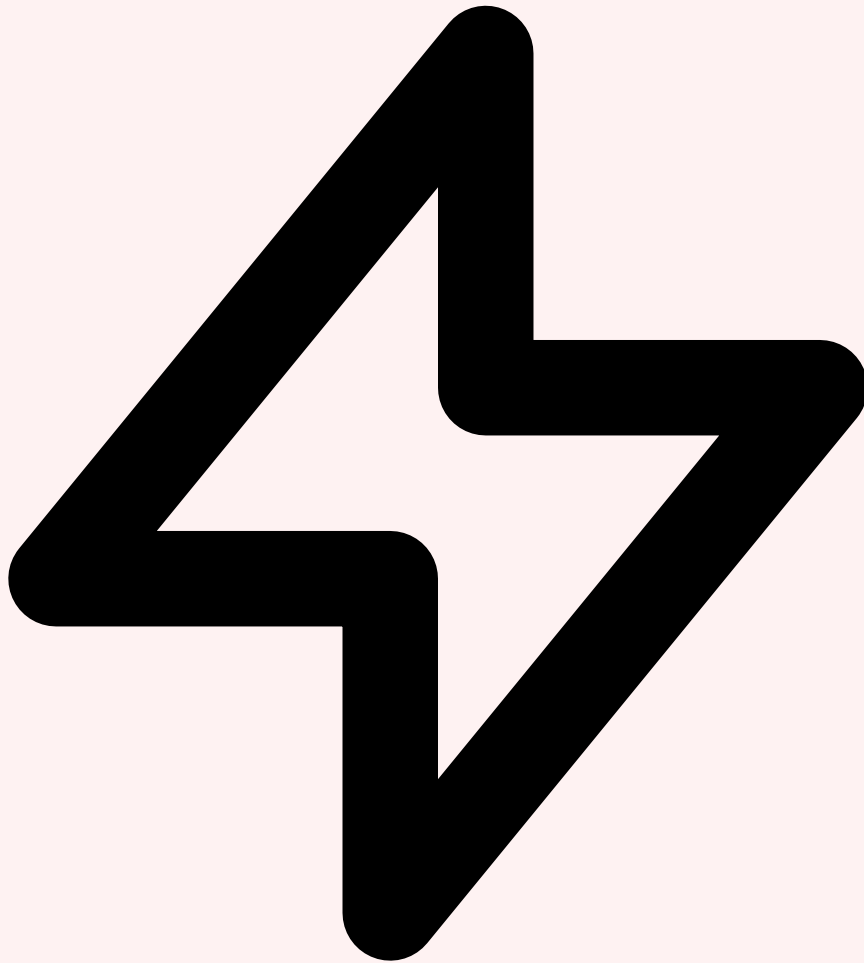


Qwen & DeepSeek Benchmarks Comparatifs **Guide de Choix**



7 Guide de Choix par Cas d'Usage

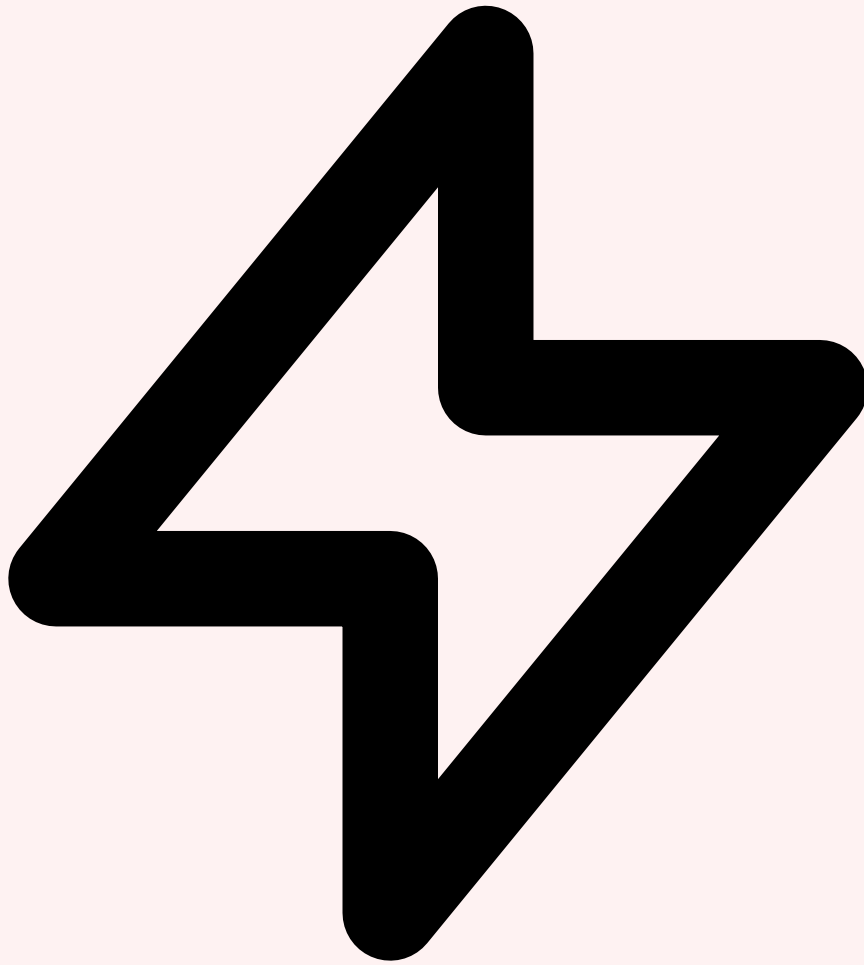
Au-delà des benchmarks, le choix d'un LLM open source doit être guidé par vos **contraintes opérationnelles concrètes** : budget matériel, cas d'usage principal, exigences réglementaires, compétences internes et besoins de personnalisation. Voici un arbre de décision pragmatique pour orienter votre choix.



Arbre de décision par budget matériel

Le premier critère de sélection est souvent le **budget matériel disponible**. Voici les recommandations par tranche d'équipement : Pour approfondir, consultez [Agents IA pour le SOC : Triage Automatisé des Alertes](#).

- **Smartphone / Raspberry Pi (2-4 Go RAM)** — Gemma 3 1B ou 4B quantifié. Seul choix viable pour l'edge computing avec des performances acceptables sur les tâches simples de classification, résumé court et Q&A basique.
- **GPU consommateur 16-24 Go (RTX 4090, RTX 5090)** — Gemma 3 27B Q4, Qwen 2.5 32B Q4, ou Mistral Small. Excellent pour le développement, le prototypage et les charges de travail modérées en production.
- **GPU professionnel 48-80 Go (A6000, H100)** — Llama 4 Scout Q4, Qwen 2.5 72B Q8, Mistral Large 2 Q4. Permet de déployer des modèles de classe professionnelle avec des performances proches du full-precision.
- **Cluster multi-GPU (2-8x H100)** — Llama 4 Maverick, DeepSeek V3 en full-precision. Pour les organisations qui exigent les meilleures performances absolues et disposent de l'infrastructure correspondante.

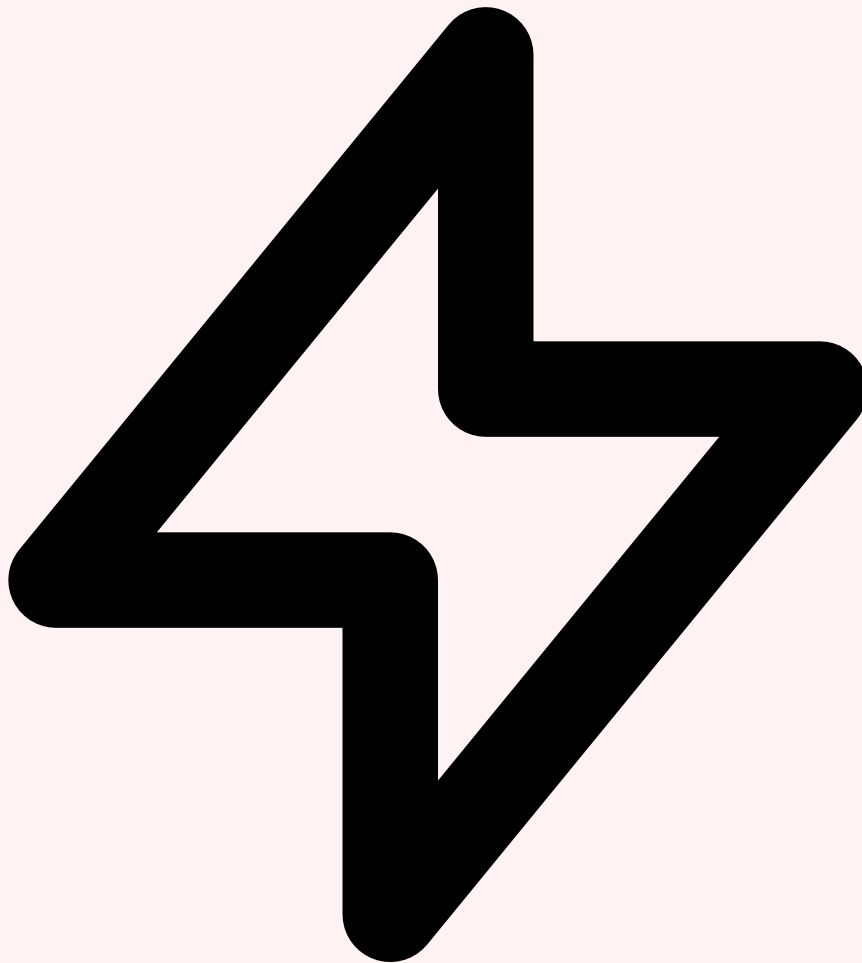


Choix par cas d'usage métier

Le cas d'usage détermine souvent le modèle optimal bien plus que les benchmarks génériques :

- **» Chatbot d'entreprise / Support client** — Llama 4 Scout pour la meilleure qualité conversationnelle, ou Qwen 2.5 72B pour un excellent rapport qualité/coût. La fenêtre de 10M tokens de Scout est un atout pour maintenir le contexte sur de longues conversations.
- **» Génération de code / Assistant développeur** — DeepSeek V3 ou Codestral de Mistral. DeepSeek excelle sur les problèmes algorithmiques complexes, Codestral sur le code idiomatique et la complétion contextuelle dans les IDE.
- **» RAG / Analyse documentaire** — Llama 4 Scout (contexte 10M tokens) ou Llama 4 Maverick pour les corpus massifs. La capacité à ingérer des documents entiers sans chunking simplifie considérablement l'architecture RAG.
- **» Contenu francophone / Conformité RGPD** — Mistral Large 2, sans hésitation. Sa maîtrise du français est inégalée et l'entreprise est soumise à la réglementation européenne, offrant une chaîne de confiance complète.

- **›Licence et liberté d'usage** — Qwen 2.5 (Apache 2.0) et Gemma 3 (licence permissive) offrent la plus grande liberté. Llama 4 impose une limite à 700M d'utilisateurs. Mistral Large 2 requiert une licence commerciale pour un usage en production.
- **›Sécurité et guardrails** — Gemma 3 avec ShieldGemma est le choix le plus robuste pour les environnements nécessitant des filtres de sécurité intégrés. Llama Guard complète Llama 4 sur ce terrain.



Recommandation finale

Si nous devons résumer nos recommandations en une seule phrase par profil d'utilisateur :

- **›Startup / PME avec budget limité** — Commencez avec Qwen 2.5 72B sous licence Apache 2.0, le meilleur rapport qualité/coût/liberté du marché.
- **›ETI / Grand groupe européen** — Adoptez Mistral Large 2 pour la conformité RGPD et l'excellence en français, complété par Llama 4 Scout pour les tâches à contexte long.

- **Équipe R&D / Recherche** — Explorez DeepSeek V3 pour sa performance brute et ses innovations architecturales. Son rapport performance/coût d'entraînement est majeur.
- **Déploiement edge / IoT** — Gemma 3 est votre seul choix viable, et c'est un excellent choix. La gamme 1B-27B couvre tous les scénarios embarqués.

Le paysage LLM open source évolue à un rythme effréné. Ce comparatif reflète l'état de l'art en février 2026, mais de nouvelles releases sont attendues chaque trimestre. Nous recommandons de réévaluer votre choix tous les six mois et de maintenir une **architecture modulaire** qui vous permet de remplacer le modèle sous-jacent sans refondre l'ensemble de votre pipeline applicatif. Les outils comme vLLM, Ollama et LiteLLM facilitent cette portabilité en fournissant une interface d'API unifiée indépendante du modèle utilisé.

Besoin d'un accompagnement expert ?

Nos consultants en cybersécurité et IA vous accompagnent dans vos projets. Devis personnalisé sous 24h.

Références et ressources externes

- OWASP LLM Top 10 — Les 10 risques majeurs pour les applications LLM
- MITRE ATLAS — Framework de menaces pour les systèmes d'intelligence artificielle
- NIST AI RMF — AI Risk Management Framework du NIST
- arXiv — Archive ouverte de publications scientifiques en IA
- HuggingFace Docs — Documentation de référence pour les modèles de ML

Pour approfondir ce sujet, consultez notre outil open-source ml-model-security-audit qui facilite l'évaluation de la sécurité des modèles ML.

Sources et références : [ArXiv IA](#) · [Hugging Face Papers](#)

FAQ

Qu'est-ce que Llama 4, Mistral Large, Gemma 3 ?

Le concept de Llama 4, Mistral Large, Gemma 3 est détaillé dans les premières sections de cet article, qui couvrent les fondamentaux, les enjeux et le contexte opérationnel. Pour un accompagnement sur ce sujet, [contactez nos experts](#).

Pourquoi Llama 4, Mistral Large, Gemma 3 est-il important en cybersécurité ?

La compréhension de Llama 4, Mistral Large, Gemma 3 permet aux équipes de sécurité d'améliorer leur posture défensive. Les sections « Table des Matières » et « 2 Llama 4 : Scout et Maverick par Meta » détaillent les raisons de cette importance. Pour un accompagnement sur ce sujet, [contactez nos experts](#).

Comment mettre en œuvre les recommandations de cet article ?

Les recommandations pratiques sont détaillées tout au long de l'article, avec des commandes, des outils et des méthodologies éprouvées. La section « Conclusion » fournit une synthèse actionnable. Pour un accompagnement sur ce sujet, [contactez nos experts](#).

Conclusion

Cet article a couvert les aspects essentiels de Table des Matières, 1 Le Paysage LLM Open Source en 2026, 2 Llama 4 : Scout et Maverick par Meta. La mise en pratique de ces recommandations permet de renforcer significativement la posture de sécurité de votre organisation.

Ayi NEDJIMI Consultants — Expert cybersécurité offensive & intelligence artificielle

ayinedjimi-consultants.fr · ayi@ayinedjimi-consultants.fr

© 2026 — Reproduction interdite sans autorisation.