

Veille IA – Mars 2025 : Nouvelles avancées et tendances majeures

Synthèse globale des avancées en IA en mars 2025

Le mois de mars 2025 a été marqué par une accélération des innovations en intelligence artificielle, tant du côté des modèles de langage géants (LLM) que des IA visuelles et robotiques. Les principaux acteurs (OpenAI, Google, Anthropic, Meta et autres) ont dévoilé des évolutions importantes de leurs modèles : OpenAI a lancé **GPT-4.5**, une version intermédiaire plus « humaine » de son modèle de conversation phare, tandis que Google DeepMind a présenté **Gemini 2.5**, un modèle « raisonneur » franchissant un nouveau palier de performance ([Gemini 2.5: Our newest Gemini model with thinking](https://blog.google/technology/google-deepmind/gemini-model-thinking-updates-march-2025/#:~:text=Today%20we%E2%80%99re%20introducing%20Gemini%202.5,LMarena%20by%20a%20significant%20margin) (<https://blog.google/technology/google-deepmind/gemini-model-thinking-updates-march-2025/#:~:text=Today%20we%E2%80%99re%20introducing%20Gemini%202.5,LMarena%20by%20a%20significant%20margin>)) (Gemini 2.5: Our newest Gemini model with thinking (<https://blog.google/technology/google-deepmind/gemini-model-thinking-updates-march-2025/#:~:text=Introducing%20Gemini%202.5>)). Dans le même temps, les modèles open source continuent de gagner du terrain : la startup française Mistral a publié un modèle léger **Small 3.1** surpassant ses concurrents de même taille ([Mistral AI lance Small 3.1, un modèle léger qui prétend surpasser la concurrence](https://www.blogdumoderateur.com/mistral-ai-lance-small-3-1-surpasser-concurrence/#:~:text=Cette%20nouvelle%20version%2C%20qui%20s%E2%80%99appuie,offrant%20une%20vitesse%20d%E2%80%99inf%C3%A9rence%20inf%C3%A9rieur) (<https://www.blogdumoderateur.com/mistral-ai-lance-small-3-1-surpasser-concurrence/#:~:text=Cette%20nouvelle%20version%2C%20qui%20s%E2%80%99appuie,offrant%20une%20vitesse%20d%E2%80%99inf%C3%A9rence%20inf%C3%A9rieur>)) et Meta a rendu disponible **Llama 3.1** (405 milliards de paramètres), le plus grand modèle ouvert à ce jour ([Introducing Llama 3.1: Our most capable models to date - Meta AI](https://ai.meta.com/blog/meta-llama-3-1/#:~:text=AI%20ai,capable%20openly%20available%20foundation%20model) (<https://ai.meta.com/blog/meta-llama-3-1/#:~:text=AI%20ai,capable%20openly%20available%20foundation%20model>)).

Les **outils et plateformes** évoluent pour faciliter l'intégration de ces IA : des frameworks comme LangChain introduisent de nouvelles bibliothèques pour construire des **agents IA autonomes** collaboratifs, avec mémoire à long terme et navigation web intégrée ([LangChain - Changelog](https://changelog.langchain.com/?page=2#:~:text=LangGraph%20LangGraph%20Swarm%20for%20building,March%201%2C%202025) (<https://changelog.langchain.com/?page=2#:~:text=LangGraph%20LangGraph%20Swarm%20for%20building,March%201%2C%202025>)). Parallèlement, l'IA générative visuelle franchit de nouveaux caps : Midjourney a déployé sa version 7 pour des images encore plus réalistes, Runway a lancé **Gen-4** pour la génération vidéo cohérente ([Runway releases an impressive new video-generating AI model](https://techcrunch.com/2025/03/31/runway-releases-an-impressive-new-video-generating-ai-model/#:~:text=highest) | TechCrunch (<https://techcrunch.com/2025/03/31/runway-releases-an-impressive-new-video-generating-ai-model/#:~:text=highest>)), et Stability AI innove en transformant de simples photos 2D en scènes 3D animées ([Stability AI's new AI model turns photos into 3D scenes](https://techcrunch.com/2025/03/18/stability-ai-new-ai-model-turns-photos-into-3d-scenes/#:~:text=Stability%20AI%20has%20released%20a,with%20realistic%20depth%20and%20perspective) | TechCrunch (<https://techcrunch.com/2025/03/18/stability-ai-new-ai-model-turns-photos-into-3d-scenes/#:~:text=Stability%20AI%20has%20released%20a,with%20realistic%20depth%20and%20perspective>)). Ces avancées techniques s'accompagnent de premières applications concrètes et de démonstrations frappantes : par exemple, en Chine, l'agent conversationnel **Manus** a impressionné par son autonomie complète dans la prise de décision ([IA : 3 grandes avancées en mars 2025 qui vont tout changer - Upskilling](https://upskilling.com/ia-3-grandes-avancees-en-mars-2025-qui-vont-tout-changer/#:~:text=Le%206%20mars%202025%2C%20la,besoin%20d%E2%80%99une%20intervention%20humaine%20constante) (<https://upskilling.com/ia-3-grandes-avancees-en-mars-2025-qui-vont-tout-changer/#:~:text=Le%206%20mars%202025%2C%20la,besoin%20d%E2%80%99une%20intervention%20humaine%20constante>)). Face à ces progrès rapides, les gouvernements et industries s'organisent : un sommet mondial à Paris a abouti à une déclaration commune pour une IA **éthique et transparente**, tandis qu'un code de bonnes pratiques a été dévoilé pour encadrer les IA généralistes ([Intelligence Artificielle : Les 5 actualités majeures du 14 mars 2025](https://fr.linkedin.com/pulse/intelligence-artificielle-les-5-actualites-majeures-du-14-mars-2025) (<https://fr.linkedin.com/pulse/intelligence-artificielle-les-5-actualites-majeures-du-14-mars-2025>)).

Les entreprises investissent massivement (OpenAI a consacré 50 M\$ à un consortium académique NextGenAI ([March 2025 round-up of interesting AI news and announcements - Artificial intelligence](https://nationalcentreforai.jiscinvolvement.org/wp/2025/03/27/march-2025-round-up-of-interesting-ai-news-and-announcements/#:~:text=Introducing%20NextGenAI%3A%20a%20consortium%20to,using%20AI%20to%20accelerate%20research) ([http://nationalcentreforai.jiscinvolvement.org/wp/2025/03/27/march-2025-round-up-of-interesting-ai-news-and-announcements/#:~:text=Introducing%20NextGenAI%3A%20a%20consortium%20to,using%20AI%20to%20accelerate%20research](https://nationalcentreforai.jiscinvolvement.org/wp/2025/03/27/march-2025-round-up-of-interesting-ai-news-and-announcements/#:~:text=Introducing%20NextGenAI%3A%20a%20consortium%20to,using%20AI%20to%20accelerate%20research)))) et s'allient pour sécuriser et fiabiliser l'IA (partenariat HuggingFace-JFrog pour analyser la **sécurité des modèles** sur le hub open source ([JFrog s'associe à Hugging Face pour assurer la sécurité des modèles GenAI](https://www.solutions-numeriques.com/jfrog-sassocie-a-hugging-face-pour-assurer-la-securite-des-modeles-genai/#:~:text=Cette%20int%C3%A9gration%20consistera%2C%20sch%C3%A9matiquement%2C%20%C3%A0,Karas%2C%20CTO%20de%20JFrog%20Security) (<https://www.solutions-numeriques.com/jfrog-sassocie-a-hugging-face-pour-assurer-la-securite-des-modeles-genai/#:~:text=Cette%20int%C3%A9gration%20consistera%2C%20sch%C3%A9matiquement%2C%20%C3%A0,Karas%2C%20CTO%20de%20JFrog%20Security>))). En somme, mars 2025 illustre la convergence des efforts : rendre l'IA plus **puissante**, plus **accessible** et plus **responsable**, tout en anticipant les défis techniques (coûts de calcul, hallucinations) et sociétaux (emploi, régulation) à venir.

Principales avancées techniques en mars 2025

□ Nouvelles générations de modèles de langage (LLM) et NLP

- **OpenAI – GPT-4.5 (2025)** : OpenAI a introduit GPT-4.5, une version améliorée de GPT-4 présentée comme « *le modèle de conversation le plus avancé à ce jour* ». Disponible en avant-première pour les abonnés ChatGPT Plus depuis fin février ([Les 10 modèles d'IA les plus performants en mars 2025](https://www.blogdumoderateur.com/modeles-ia-plus-performants-mars-2025/#:~:text=conserve%20sa%20premi%C3%A8re%20position%20acquise,cinq%20ans%20plus%20t%C3%A0,qu%E2%80%99a%20destin%C3%A9) (<https://www.blogdumoderateur.com/modeles-ia-plus-performants-mars-2025/#:~:text=conserve%20sa%20premi%C3%A8re%20position%20acquise,cinq%20ans%20plus%20t%C3%A0,qu%E2%80%99a%20destin%C3%A9>)) (Introducing GPT-4.5 | OpenAI (<https://openai.com/index/introducing-gpt-4-5/#:~:text=We%E2%80%99re%20releasing%20a%20research%20preview,generate%20creative%20insights%20without%20reasoning>))) , GPT-4.5 apporte des réponses plus **naturelles et nuancées**, avec une meilleure compréhension des intentions et émotions humaine ([OpenAI dévoile GPT-4.5, plus intelligent et plus "humain"](https://www.mac4ever.com/ia/187546-openai-devoile-gpt-4-5-plus-intelligent-et-plus-humain/#:~:text=Avec%20GPT,y%20compris%20les%20attentes%20implicites) (<https://www.mac4ever.com/ia/187546-openai-devoile-gpt-4-5-plus-intelligent-et-plus-humain/#:~:text=Avec%20GPT,y%20compris%20les%20attentes%20implicites>))) . Par exemple, face à un utilisateur frustré, GPT-4.5 sait formuler une réponse empathique plutôt que littérale, témoignant de son « *intelligence émotionnelle* » accru ([OpenAI dévoile GPT-4.5, plus intelligent et plus "humain"](https://www.mac4ever.com/ia/187546-openai-devoile-gpt-4-5-plus-intelligent-et-plus-humain/#:~:text=Un%20exemple%20frappant%20mis%20en,%20compte%20le%20contexte%20%C3%A9motionnel) (<https://www.mac4ever.com/ia/187546-openai-devoile-gpt-4-5-plus-intelligent-et-plus-humain/#:~:text=Un%20exemple%20frappant%20mis%20en,%20compte%20le%20contexte%20%C3%A9motionnel>))) . Conçu comme étape intermédiaire avant GPT-5, ce modèle plus grand réduit les hallucinations et améliore le suivi des instructions ([Introducing GPT-4.5 | OpenAI](https://www.inc.com/ben-sherry/openai-shifts-course-says-gpt-5-coming-in-a-few-months) (<https://openai.com/index/introducing-gpt-4-5/#:~:text=We%E2%80%99re%20releasing%20a%20research%20preview,generate%20creative%20insights%20without%20reasoning>))) . (À noter : OpenAI a indiqué que GPT-5, qui intégrera pleinement les capacités de raisonnement « O », est prévu d'ici quelques mois ([OpenAI Shifts Course, Says GPT-5 Coming in 'a Few Months'](https://www.inc.com/ben-sherry/openai-shifts-course-says-gpt-5-coming-in-a-few-months) (<https://www.inc.com/ben-sherry/openai-shifts-course-says-gpt-5-coming-in-a-few-months>)).

[months/91172067#:~:text=rolling%20out%20its%20next%20AI,5%2C%20in%20%E2%80%9Ca%20few%20months.%E2%80%9D\)\)](#) (OpenAI Shifts Course, Says GPT-5 Coming in 'a Few Months' ([https://www.inc.com/ben-sherry/openai-shifts-course-says-gpt-5-coming-in-a-few-months/91172067#:~:text=The%20capacity%20concerns%20come%20on.%E2%80%9D\)\)](https://www.inc.com/ben-sherry/openai-shifts-course-says-gpt-5-coming-in-a-few-months/91172067#:~:text=The%20capacity%20concerns%20come%20on.%E2%80%9D)))) .

- **OpenAI – modèles “O” et raisonnement** : En parallèle, OpenAI fait évoluer sa gamme de modèles dédiés au **raisonnement pas-à-pas** (“O” pour *Orion*). Initialement, l'idée était de fusionner ces modèles dans GPT-5, mais l'entreprise a finalement décidé de les publier séparément pour affiner la transitio ([OpenAI Shifts Course, Says GPT-5 Coming in 'a Few Months'](#) ([https://www.inc.com/ben-sherry/openai-shifts-course-says-gpt-5-coming-in-a-few-months/91172067#:~:text=In%20February%2C%20Altman%20posted%20on.5\)\)](https://www.inc.com/ben-sherry/openai-shifts-course-says-gpt-5-coming-in-a-few-months/91172067#:~:text=In%20February%2C%20Altman%20posted%20on.5)))) . Ainsi, le modèle **o3** sera lancé comme agent de raisonnement avancé, accompagné de **o4-mini** (version moins coûteuse ([OpenAI Shifts Course, Says GPT-5 Coming in 'a Few Months'](#) ([https://www.inc.com/ben-sherry/openai-shifts-course-says-gpt-5-coming-in-a-few-months/91172067#:~:text=a%20standalone%20product%20and%20would.5\)\)](https://www.inc.com/ben-sherry/openai-shifts-course-says-gpt-5-coming-in-a-few-months/91172067#:~:text=a%20standalone%20product%20and%20would.5)))) . L'objectif est de doter les IA conversationnelles d'une capacité de **chaînage logique** – par exemple résoudre des problèmes de maths complexes en décomposant les étapes. OpenAI intègre déjà certaines de ces capacités à ChatGPT : on a vu apparaître un mode “**Deep Research**” capable d'enquêter de manière autonome sur un sujet point ([OpenAI's Latest Breakthrough: AI That Comes Up With New Ideas](#) ([https://www.theinformation.com/articles/openai-latest-breakthrough-ai-comes-new-ideas#:~:text=OpenAI%27s%20latest%20Breakthrough%3A%20AI%20That,as%20deep%20research%2C%20which\)\)](https://www.theinformation.com/articles/openai-latest-breakthrough-ai-comes-new-ideas#:~:text=OpenAI%27s%20latest%20Breakthrough%3A%20AI%20That,as%20deep%20research%2C%20which)))) . Ces efforts vers des *AI agents* plus autonomes s'inscrivent dans une tendance forte du mois : la montée de ce qu'OpenAI et d'autres appellent l’**“IA agentique”**, c'est-à-dire des IA pouvant planifier et agir de leur propre initiative pour accomplir des objectifs complexes.
- **Google DeepMind – Gemini 2.5** : Google a frappé fort en dévoilant Gemini 2.5 Pro, décrit comme « *notre modèle d'IA le plus intelligent* ». Contrairement aux versions précédentes de Gemini, axées sur la génération, la série 2.x inaugure les “*thinking models*” intégrant nativement un mécanisme de **raisonnement interne** (inspiration proche du *chain-of-thought*). Le modèle Gemini 2.5 Pro (version expérimentale) a pris la **1ère place** du classement LMArena/Chatbot Arena en mars, devançant GPT-4.5 et consort ([Gemini 2.5: Our newest Gemini model with thinking](#) ([https://blog.google/technology/google-deepmind/gemini-model-thinking-updates-march-2025/#:~:text=Today%20we%E2%80%99re%20introducing%20Gemini%202,LMArena%20by%20a%20significant%20margin\)\)](https://blog.google/technology/google-deepmind/gemini-model-thinking-updates-march-2025/#:~:text=Today%20we%E2%80%99re%20introducing%20Gemini%202,LMArena%20by%20a%20significant%20margin)))) ([Gemini 2.5: Our newest Gemini model with thinking](#) ([https://blog.google/technology/google-deepmind/gemini-model-thinking-updates-march-2025/#:~:text=Introducing%20Gemini%202\)\)](https://blog.google/technology/google-deepmind/gemini-model-thinking-updates-march-2025/#:~:text=Introducing%20Gemini%202)))) . Il excelle sur les tâches complexes : **codage**, mathématiques et questions scientifiques pointues, où il dépasse GPT-4.5 et Claude 3.7 dans plusieurs benchmark ([Gemini 2.5: Our newest Gemini model with thinking](#) ([https://blog.google/technology/google-deepmind/gemini-model-thinking-updates-march-2025/#:~:text=to%20capture%20the%20human%20frontier.of%20knowledge%20and%20reasoning\)\)](https://blog.google/technology/google-deepmind/gemini-model-thinking-updates-march-2025/#:~:text=to%20capture%20the%20human%20frontier.of%20knowledge%20and%20reasoning)))) ([Gemini 2.5: Our newest Gemini model with thinking](#) ([https://blog.google/technology/google-deepmind/gemini-model-thinking-updates-march-2025/#:~:text=Image%3A%20Bar%20charts%20comparing%20the,strong%20results%20in%20all%20categories\)\)](https://blog.google/technology/google-deepmind/gemini-model-thinking-updates-march-2025/#:~:text=Image%3A%20Bar%20charts%20comparing%20the,strong%20results%20in%20all%20categories)))) . Disponible via Google AI Studio pour les clients *Gemini Advanced*, ce modèle reflète la stratégie de Google : combiner puissance brute et réflexion structurée. En pratique, cela se traduit par une meilleure cohérence dans les réponses nécessitant analyse et contexte. (Google a aussi indiqué que ces capacités de “*thinking*” seront généralisées à tous ses futurs modèles, grands et petit ([Gemini 2.5: Our newest Gemini model with thinking](#) ([https://blog.google/technology/google-deepmind/gemini-model-thinking-updates-march-2025/#:~:text=For%20a%20long%20time%2C%20we%E2%80%99ve.0%20Flash%20Thinking\)\)](https://blog.google/technology/google-deepmind/gemini-model-thinking-updates-march-2025/#:~:text=For%20a%20long%20time%2C%20we%E2%80%99ve.0%20Flash%20Thinking)))) .) Par ailleurs, la marque **Bard** a laissé place au nom *Gemini* dans ses produits grand public : l'assistant AI de Google utilise désormais les modèles Gemini pour offrir des réponses plus fiables et multimodale ([Gemini - Google](#) ([https://gemini.google.com/#:~:text=Bard%20is%20now%20Gemini,and%20more%20from%20Google%20AI\)\)](https://gemini.google.com/#:~:text=Bard%20is%20now%20Gemini,and%20more%20from%20Google%20AI)))) .
- **Anthropic – Claude (évolutions)** : Si Anthropic n'a pas lancé de *Claude 3* en mars, il a notamment enrichi son assistant **Claude 2/3** existant avec de **nouvelles capacités**. La plus marquante est l'ajout du **recherche web en direct** pour Claude 3.7 (*Claude 3.7 Sonnet*), comblant un retard vis-à-vis de ChatGPT. Depuis le 20 mars, les utilisateurs payants de Claude peuvent activer un mode navigation : l'IA effectue alors des recherches internet et fournit des réponses à jour avec des “*sources citées*” ([Anthropic adds web search to its Claude chatbot | TechCrunch](#) ([https://techcrunch.com/2025/03/20/anthropic-adds-web-search-to-its-claude-chatbot/#:~:text=that%20had%20long%20eluded%20it\)\)](https://techcrunch.com/2025/03/20/anthropic-adds-web-search-to-its-claude-chatbot/#:~:text=that%20had%20long%20eluded%20it)))) ([Anthropic adds web search to its Claude chatbot | TechCrunch](#) ([https://techcrunch.com/2025/03/20/anthropic-adds-web-search-to-its-claude-chatbot/#:~:text=%E2%80%9CWhen%20Claude%20incorporates%20information%20from.%E2%80%9D\)\)](https://techcrunch.com/2025/03/20/anthropic-adds-web-search-to-its-claude-chatbot/#:~:text=%E2%80%9CWhen%20Claude%20incorporates%20information%20from.%E2%80%9D)))) . Cela permet à Claude de sortir de sa base de connaissances statique et de s'appuyer sur des informations en temps réel, avec des références que l'utilisateur peut vérifier. Par exemple, interrogé sur un événement tech de la semaine, Claude va chercher les actualités pertinentes (sur X/Twitter, sites d'actus, etc.) et produire une synthèse sourcée (voir l'illustration ci-dessous). Ce **rattrapage fonctionnel** met Claude au niveau de ChatGPT/Bing et Google Gemini sur ce point ([Anthropic adds web search to its Claude chatbot | TechCrunch](#) ([https://techcrunch.com/2025/03/20/anthropic-adds-web-search-to-its-claude-chatbot/#:~:text=Claude%E2%80%99s%20ability%20to%20search%20the,with%20the%20reversal%20in%20course\)\)](https://techcrunch.com/2025/03/20/anthropic-adds-web-search-to-its-claude-chatbot/#:~:text=Claude%E2%80%99s%20ability%20to%20search%20the,with%20the%20reversal%20in%20course)))) . Anthropic a également intégré Claude à des outils de travail : en mars, **Claude Pro** s'est connecté à Gmail, Google Docs et Calendar pour aider à synthétiser emails et documents professionnel ([Claude takes research to new places \ Anthropic](#) ([https://www.anthropic.com/news/research#:~:text=Google%20Workspace\)\)](https://www.anthropic.com/news/research#:~:text=Google%20Workspace)))) ([Claude takes research to new places \ Anthropic](#) ([https://www.anthropic.com/news/research#:~:text=Claude%20now%20integrates%20with%20Gmail,about%20your%20work%20and%20schedule\)\)](https://www.anthropic.com/news/research#:~:text=Claude%20now%20integrates%20with%20Gmail,about%20your%20work%20and%20schedule)))) . Cet **assistant augmenté** peut, par exemple, résumer vos derniers échanges mail et rechercher des infos en ligne pour préparer un meeting – signe que les IA deviennent de plus en plus **imbriquées dans les flux de travail**.
([Anthropic adds web search to its Claude chatbot | TechCrunch](#) (<https://techcrunch.com/2025/03/20/anthropic-adds-web-search-to-its-claude-chatbot/>)) **Exemple** : interface de **Claude 3.7** effectuant une recherche web pour répondre à une question d'actualité (résultats avec sources affichées ([Anthropic adds web search to its Claude chatbot | TechCrunch](#) ([https://techcrunch.com/2025/03/20/anthropic-adds-web-search-to-its-claude-chatbot/#:~:text=%E2%80%9CWhen%20Claude%20incorporates%20information%20from.%E2%80%9D\)\)](https://techcrunch.com/2025/03/20/anthropic-adds-web-search-to-its-claude-chatbot/#:~:text=%E2%80%9CWhen%20Claude%20incorporates%20information%20from.%E2%80%9D)))) ([Anthropic adds web search to its Claude chatbot | TechCrunch](#) ([https://techcrunch.com/2025/03/20/anthropic-adds-web-search-to-its-claude-chatbot/#:~:text=Claude%E2%80%99s%20ability%20to%20search%20the,with%20the%20reversal%20in%20course\)\)](https://techcrunch.com/2025/03/20/anthropic-adds-web-search-to-its-claude-chatbot/#:~:text=Claude%E2%80%99s%20ability%20to%20search%20the,with%20the%20reversal%20in%20course))))) .
- **xAI – Grok 3** : L'initiative d'Elon Musk, moins médiatisée que les précédentes, a néanmoins fait parler d'elle en mars : **Grok-3** (modèle d'xAI accessible sur la plateforme X/Twitter) a brièvement occupé la **1ère place** de la Chatbot Arena début mar ([Les 10 modèles d'IA les plus performants en mars 2025](#) ([https://www.blogdumoderateur.com/modeles-ia-plus-performants-mars-2025/#:~:text=Malgr%C3%A9%20des%20lacunes%20apparentes%20lors,4.5%20est%20l%E2%80%99un%20des%20trois\)\)](https://www.blogdumoderateur.com/modeles-ia-plus-performants-mars-2025/#:~:text=Malgr%C3%A9%20des%20lacunes%20apparentes%20lors,4.5%20est%20l%E2%80%99un%20des%20trois)))) . Bien que la démonstration initiale ait montré des limites et quelques lacunes, Grok 3 bénéficie d'une forte exposition auprès du grand public via X et illustre l'arrivée de nouveaux acteurs sur le terrain des LLM. Sa spécialité mise en avant est la **gratuité** d'accès (avec certaines limites) et un entraînement orienté sur l'actualité immédiate – Musk ayant affirmé vouloir un chatbot « *rebelle et à jour* » pour rivaliser avec ChatGPT. Si Grok reste en retrait face à GPT-4.5 ou Gemini en performance brute, sa progression rapide en quelques mois (de Grok-1 en fin 2024 à Grok-3 en 02/2025) montre la féroce compétition en cours et la **démocratisation** des grands modèles.

- Meta – LLaMA 3.1 (open source)** : Du côté des modèles ouverts, Meta a continué sur sa lancée de l'open science. Après LLaMA 2 (2023), la société a mis à disposition début 2025 **LLaMA 3** et sa version affûtée **LLaMA 3.1**. Le modèle haut de gamme Llama 3.1 affiche **405 milliards** de paramètres – ce qui en fait le plus grand modèle open source disponible publiquement ([Introducing Llama 3.1: Our most capable models to date - Meta AI](https://ai.meta.com/blog/meta-llama-3-1/#:~:text=AI%20ai%20openly%20available%20foundation%20model))) ([1#:#:~:text=AI%20ai%20openly%20available%20foundation%20model\)\)](https://ai.meta.com/blog/meta-llama-3-1/#:~:text=AI%20ai%20openly%20available%20foundation%20model)))) – et a été entraîné sur un contexte étendu (jusqu'à 128k tokens). Des premières évaluations montrent qu'il rivalise avec les meilleurs modèles fermés sur de nombreux benchmarks, tout en étant **auto-hébergeable** et modifiable à volonté ([Meta releases new Llama 3.1 models, including highly anticipated 405B parameter variant | IBM](https://www.ibm.com/think/news/meta-releases-llama-3-1-models-405b-parameter-variant#:~:text=Gemini%20Ultra%201.0%20%2853.2))) ([https://www.ibm.com/think/news/meta-releases-llama-3-1-models-405b-parameter-variant#:~:text=Gemini%20Ultra%201.0%20%2853.2\)\)](https://www.ibm.com/think/news/meta-releases-llama-3-1-models-405b-parameter-variant#:~:text=Gemini%20Ultra%201.0%20%2853.2)))) ([Meta releases new Llama 3.1 models, including highly anticipated 405B parameter variant | IBM](https://www.ibm.com/think/news/meta-releases-llama-3-1-models-405b-parameter-variant#:~:text=Gemini%20Ultra%201.0%20%2853.2))) ([https://www.ibm.com/think/news/meta-releases-llama-3-1-models-405b-parameter-variant#:~:text=Gemini%20Ultra%201.0%20%2853.2\)\)](https://www.ibm.com/think/news/meta-releases-llama-3-1-models-405b-parameter-variant#:~:text=Gemini%20Ultra%201.0%20%2853.2))))]. Par exemple, son score en Q&A avancé (jeu de questions GPQA) atteint ~50.7%, au niveau de Claude 3 (50.4%) et GPT-4 "T" ([Meta releases new Llama 3.1 models, including highly anticipated 405B parameter variant | IBM](https://www.ibm.com/think/news/meta-releases-llama-3-1-models-405b-parameter-variant#:~:text=Gemini%20Ultra%201.0%20%2853.2))) ([https://www.ibm.com/think/news/meta-releases-llama-3-1-models-405b-parameter-variant#:~:text=Gemini%20Ultra%201.0%20%2853.2\)\)](https://www.ibm.com/think/news/meta-releases-llama-3-1-models-405b-parameter-variant#:~:text=Gemini%20Ultra%201.0%20%2853.2))))]. Sa disponibilité sur des plateformes comme IBM Watsonx ou Oracle OCI a été annoncée courant mar ([Meta releases new Llama 3.1 models, including highly anticipated 405B parameter variant | IBM](https://www.ibm.com/think/news/meta-releases-llama-3-1-models-405b-parameter-variant#:~:text=Gemini%20Ultra%201.0%20%2853.2))) ([https://www.ibm.com/think/news/meta-releases-llama-3-1-models-405b-parameter-variant#:~:text=Gemini%20Ultra%201.0%20%2853.2\)\)](https://www.ibm.com/think/news/meta-releases-llama-3-1-models-405b-parameter-variant#:~:text=Gemini%20Ultra%201.0%20%2853.2))))], signe d'une adoption industrielle. Llama 3.1 est appelé une *foundation model* "stable" : les développeurs peuvent l'utiliser comme base et le fine-tuner selon leurs besoins spécifiques, sans craindre des changements arbitraires d'AP ([Meta releases new Llama 3.1 models, including highly anticipated 405B parameter variant | IBM](https://www.ibm.com/think/news/meta-releases-llama-3-1-models-405b-parameter-variant#:~:text=Gemini%20Ultra%201.0%20%2853.2))) ([https://www.ibm.com/think/news/meta-releases-llama-3-1-models-405b-parameter-variant#:~:text=Gemini%20Ultra%201.0%20%2853.2\)\)](https://www.ibm.com/think/news/meta-releases-llama-3-1-models-405b-parameter-variant#:~:text=Gemini%20Ultra%201.0%20%2853.2))))]. Cette ouverture, opposée à la **boîte noire** des modèles SaaS, est saluée comme un atout pour la recherche et la reproductibilité scientifique ([Meta releases new Llama 3.1 models, including highly anticipated 405B parameter variant | IBM](https://www.ibm.com/think/news/meta-releases-llama-3-1-models-405b-parameter-variant#:~:text=Gemini%20Ultra%201.0%20%2853.2))) ([https://www.ibm.com/think/news/meta-releases-llama-3-1-models-405b-parameter-variant#:~:text=Gemini%20Ultra%201.0%20%2853.2\)\)](https://www.ibm.com/think/news/meta-releases-llama-3-1-models-405b-parameter-variant#:~:text=Gemini%20Ultra%201.0%20%2853.2))))]. En parallèle, Meta a intégré Llama 3 dans ses propres produits : l'assistant *Meta AI* des plateformes Facebook/Instagram repose sur ces modèles et a gagné en capacités multilingues et visuelles (génération d'images, etc.).
- Mistral AI – Small 3.1 (open source)** : Autre réussite européenne, la startup française Mistral AI a lancé le 17 mars son nouveau modèle **Mistral Small 3.1* ([Mistral AI lance Small 3.1, un modèle léger qui prétend surpasser la concurrence](https://www.blogdumoderateur.com/mistral-ai-lance-small-3-1-surpasser-concurrence/#:~:text=Mistral%20AI%20a%20annonc%C3%A9%20ce%20mod%C3%A8les%20de%20sa%20cat%C3%A9gorie))) ([https://www.blogdumoderateur.com/mistral-ai-lance-small-3-1-surpasser-concurrence/#:~:text=Mistral%20AI%20a%20annonc%C3%A9%20ce%20mod%C3%A8les%20de%20sa%20cat%C3%A9gorie\)\)](https://www.blogdumoderateur.com/mistral-ai-lance-small-3-1-surpasser-concurrence/#:~:text=Mistral%20AI%20a%20annonc%C3%A9%20ce%20mod%C3%A8les%20de%20sa%20cat%C3%A9gorie))))]. Il s'agit d'un **petit LLM** (catégorie ~7 milliards de paramètres, nommé "24B" dans les tests car équivalent 24 milliards densifiés) optimisé pour tourner en local avec des ressources modestes. Open source sous licence Apache 2.0, on peut le télécharger librement sur Hugging Face et même l'exécuter sur un PC équipé d'une seule GPU grand public (une RTX 4090 suffit ([Mistral AI lance Small 3.1, un modèle léger qui prétend surpasser la concurrence](https://www.blogdumoderateur.com/mistral-ai-lance-small-3-1-surpasser-concurrence/#:~:text=Mistral%20Small%203%20est%20%C3%A9galement%20%C2%BB))) ([https://www.blogdumoderateur.com/mistral-ai-lance-small-3-1-surpasser-concurrence/#:~:text=Mistral%20Small%203%20est%20%C3%A9galement%20%C2%BB\)\)](https://www.blogdumoderateur.com/mistral-ai-lance-small-3-1-surpasser-concurrence/#:~:text=Mistral%20Small%203%20est%20%C3%A9galement%20%C2%BB))))]. Malgré sa taille réduite, Mistral 3.1 impressionne : grâce à des optimisations d'architecture, il atteint des performances **supérieures aux autres modèles de même classe** (30–40 milliards). Par exemple, sur des questions pointues (benchmark GPQA-Diamond), il obtient un score de 44% là où GPT-4o Mini plafonne à 40% et Claude 3.5 Haiku à 37 ([Mistral AI lance Small 3.1, un modèle léger qui prétend surpasser la concurrence](https://www.blogdumoderateur.com/mistral-ai-lance-small-3-1-surpasser-concurrence/#:~:text=textuelles%20encore%20meilleures%20ainsi%20qu%E2%80%99une%20vitesse%20d%E2%80%99inf%C3%A9rence%20inf%C3%A) (<https://www.blogdumoderateur.com/mistral-ai-lance-small-3-1-surpasser-concurrence/#:~:text=textuelles%20encore%20meilleures%20ainsi%20qu%E2%80%99une%20vitesse%20d%E2%80%99inf%C3%A9rence%20inf%C3%A> Surtout, il maintient une latence d'inférence très faible (~11 ms/token) le rendant réactif ([Mistral AI lance Small 3.1, un modèle léger qui prétend surpasser la concurrence](https://www.blogdumoderateur.com/mistral-ai-lance-small-3-1-surpasser-concurrence/#:~:text=textuelles%20encore%20meilleures%20ainsi%20qu%E2%80%99une%20vitesse%20d%E2%80%99inf%C3%A9rence%20inf%C3%A) (<https://www.blogdumoderateur.com/mistral-ai-lance-small-3-1-surpasser-concurrence/#:~:text=textuelles%20encore%20meilleures%20ainsi%20qu%E2%80%99une%20vitesse%20d%E2%80%99inf%C3%A9rence%20inf%C3%A> La figure ci-dessous, publiée par Mistral, montre ce positionnement : Small 3.1 combine **haute connaissance** et **vitesse**, surpassant à la fois *GPT-4o Mini* et *Gemma-3 (27B)* de ses concurrents directs. En plus du texte, Mistral 3.1 gère en entrée des **images** (capacité multimodale) et peut exploiter de très longs contextes (fenêtre étendue à 128k tokens ([Mistral AI lance Small 3.1, un modèle léger qui prétend surpasser la concurrence](https://www.blogdumoderateur.com/mistral-ai-lance-small-3-1-surpasser-concurrence/#:~:text=Cette%20nouvelle%20version%2C%20qui%20s%E2%80%99appuie%20sur%20une%20vitesse%20d%E2%80%99inf%C3%A9rence%20inf%C3%A) (<https://www.blogdumoderateur.com/mistral-ai-lance-small-3-1-surpasser-concurrence/#:~:text=Cette%20nouvelle%20version%2C%20qui%20s%E2%80%99appuie%20sur%20une%20vitesse%20d%E2%80%99inf%C3%A9rence%20inf%C3%A> Ces atouts en font un candidat idéal pour des applications embarquées ou privées (assistants personnels, agents autonomes sur mobile, etc.), domaine visé par Mistral AI. Le modèle est disponible en deux variantes (base et instruct) et également via API et sur Google Cloud Vertex A ([Mistral AI lance Small 3.1, un modèle léger qui prétend surpasser la concurrence](https://www.blogdumoderateur.com/mistral-ai-lance-small-3-1-surpasser-concurrence/#:~:text=Il%20est%20aussi%20possible%20de%20prochaines%20semaines%2C%20A%C2%BB%2C%20annonc%C3%A9%20Mistral%20AI))) ([https://www.blogdumoderateur.com/mistral-ai-lance-small-3-1-surpasser-concurrence/#:~:text=Il%20est%20aussi%20possible%20de%20prochaines%20semaines%2C%20A%C2%BB%2C%20annonc%C3%A9%20Mistral%20AI\)\)](https://www.blogdumoderateur.com/mistral-ai-lance-small-3-1-surpasser-concurrence/#:~:text=Il%20est%20aussi%20possible%20de%20prochaines%20semaines%2C%20A%C2%BB%2C%20annonc%C3%A9%20Mistral%20AI))))]. ([Mistral AI lance Small 3.1, un modèle léger qui prétend surpasser la concurrence](https://www.blogdumoderateur.com/mistral-ai-lance-small-3-1-surpasser-concurrence/#:~:text=Il%20est%20aussi%20possible%20de%20prochaines%20semaines%2C%20A%C2%BB%2C%20annonc%C3%A9%20Mistral%20AI))) ([https://www.blogdumoderateur.com/mistral-ai-lance-small-3-1-surpasser-concurrence/#:~:text=Il%20est%20aussi%20possible%20de%20prochaines%20semaines%2C%20A%C2%BB%2C%20annonc%C3%A9%20Mistral%20AI\)\)](https://www.blogdumoderateur.com/mistral-ai-lance-small-3-1-surpasser-concurrence/#:~:text=Il%20est%20aussi%20possible%20de%20prochaines%20semaines%2C%20A%C2%BB%2C%20annonc%C3%A9%20Mistral%20AI))))]. Comparatif de performance vs latence : **Mistral Small 3.1** offre un meilleur score de connaissance (axe vertical) pour une latence similaire ou inférieure à des modèles concurrents plus grand ([Mistral AI lance Small 3.1, un modèle léger qui prétend surpasser la concurrence](https://www.blogdumoderateur.com/mistral-ai-lance-small-3-1-surpasser-concurrence/#:~:text=textuelles%20encore%20meilleures%20ainsi%20qu%E2%80%99une%20vitesse%20d%E2%80%99inf%C3%A9rence%20inf%C3%A) (<https://www.blogdumoderateur.com/mistral-ai-lance-small-3-1-surpasser-concurrence/#:~:text=textuelles%20encore%20meilleures%20ainsi%20qu%E2%80%99une%20vitesse%20d%E2%80%99inf%C3%A9rence%20inf%C3%A> ([Mistral AI lance Small 3.1, un modèle léger qui prétend surpasser la concurrence](https://www.blogdumoderateur.com/mistral-ai-lance-small-3-1-surpasser-concurrence/#:~:text=surpasse%20des%20mod%C3%A8les%20similaires%20comme%20offrant%20une%20vitesse%20d%E2%80%99inf%C3%A9rence%20inf%C3%A9rieur) (<https://www.blogdumoderateur.com/mistral-ai-lance-small-3-1-surpasser-concurrence/#:~:text=surpasse%20des%20mod%C3%A8les%20similaires%20comme%20offrant%20une%20vitesse%20d%E2%80%99inf%C3%A9rence%20inf%C3%A9rieur> (Extrait des tests internes publiés par Mistral AI)).
- Nouveaux modèles spécialisés** : Mars 2025 a également vu l'émergence de modèles IA ciblant des besoins spécifiques. En Europe, le projet **EuroBert** a été annoncé comme encodeur multilingue optimisé pour les langues européennes minoritaires ([Intelligence Artificielle : Les 5 actualités majeures du 14 mars 2025](https://fr.linkedin.com/pulse/intelligence-artificielle-les-5-actualit%C3%A9s-majeures-du-nobili-mpeve#:~:text=3%EF%B8%8F%E2%83%A3%20EuroBert%203%20un%20bond%20du%20langage%20naturel%20en%20Europe))) ([https://fr.linkedin.com/pulse/intelligence-artificielle-les-5-actualit%C3%A9s-majeures-du-nobili-mpeve#:~:text=3%EF%B8%8F%E2%83%A3%20EuroBert%203%20un%20bond%20du%20langage%20naturel%20en%20Europe\)\)](https://fr.linkedin.com/pulse/intelligence-artificielle-les-5-actualit%C3%A9s-majeures-du-nobili-mpeve#:~:text=3%EF%B8%8F%E2%83%A3%20EuroBert%203%20un%20bond%20du%20langage%20naturel%20en%20Europe))))]. Ce modèle de **NLP** vise à améliorer traduction automatique et analyse de texte pour les administrations et entreprises du continent, renforçant la **souveraineté numérique** face aux modèles anglo-centré ([Intelligence Artificielle : Les 5 actualités majeures du 14 mars 2025](https://fr.linkedin.com/pulse/intelligence-artificielle-les-5-actualit%C3%A9s-majeures-du-nobili-mpeve#:~:text=L%E2%80%99intelligence%20artificielle%20continue%20de%20s%E2%80%99imposer%20%C3%A0%20l%E2%80%99information%20dans%20toutes%20les) (<https://fr.linkedin.com/pulse/intelligence-artificielle-les-5-actualit%C3%A9s-majeures-du-nobili-mpeve#:~:text=L%E2%80%99intelligence%20artificielle%20continue%20de%20s%E2%80%99imposer%20%C3%A0%20l%E2%80%99information%20dans%20toutes%20les> Par ailleurs, Nvidia, en collaboration avec des labs de recherche, a présenté le modèle **Nemotron** dédié au **raisonnement mult-étapes on-demand** (résolution de problèmes scientifiques complexes) pour ses plateformes GP ([GTC 2025 – Announcements and Live Updates | NVIDIA Blog](https://blogs.nvidia.com/blog/nvidia-keynote-at-gtc-2025-ai-news-live-updates/?utm_source=chatgpt.com#:~:text=Nemotron%20reasoning%20family%20delivers%20on%20making))) ([https://blogs.nvidia.com/blog/nvidia-keynote-at-gtc-2025-ai-news-live-updates/?utm_source=chatgpt.com#:~:text=Nemotron%20reasoning%20family%20delivers%20on%20making\)\)](https://blogs.nvidia.com/blog/nvidia-keynote-at-gtc-2025-ai-news-live-updates/?utm_source=chatgpt.com#:~:text=Nemotron%20reasoning%20family%20delivers%20on%20making))))]. Enfin, Anthropic commercialise de nouvelles variantes de Claude 3.5 (Haiku, Sonnet, Opus) adaptées à différents usages : par exemple **Claude 3.5 Haiku** privilégie la rapidité pour l'étiquetage de données, quand **Claude 3.5 Sonnet** vise une intelligence maximale pour les tâches créatives et conversationnelle ([Claude 3.5 Haiku - Anthropic](https://www.anthropic.com/claude/haiku#:~:text=Claude%203%20extraction%20and%20automated%20labeling%20tasks))) ([https://www.anthropic.com/claude/haiku#:~:text=Claude%203%20extraction%20and%20automated%20labeling%20tasks\)\)](https://www.anthropic.com/claude/haiku#:~:text=Claude%203%20extraction%20and%20automated%20labeling%20tasks)))) (Amazon Bedrock introduit Claude 3.5 Haiku

and an upgraded ([https://www.aboutamazon.com/news/aws/amazon-bedrock-anthropic-ai-claude-3-5-sonnet#:~:text=Anthropic%27s%20Claude%203.use%20capability%20in%20public\)\)](https://www.aboutamazon.com/news/aws/amazon-bedrock-anthropic-ai-claude-3-5-sonnet#:~:text=Anthropic%27s%20Claude%203.use%20capability%20in%20public)))】. Cette diversification de l'écosystème de modèles permet de mieux répondre aux cas d'usage variés : plutôt que des LLM uniques polyvalents, on voit émerger une *galaxie de modèles spécialisés* (par taille, par langue, par fonction).

□ Progrès en IA générative visuelle, image et vidéo

- **Midjourney v7** : Plus d'un an après Midjourney v6, la célèbre IA de génération d'images a reçu une mise à jour majeure en mars 2025 avec la sortie de **Midjourney V7** ([Test de Midjourney v7 : un modèle bourré de qualités... mais aussi de défauts](https://www.blogdumoderateur.com/test-midjourney-v7#:~:text=styles%20graphiques))) ([https://www.blogdumoderateur.com/test-midjourney-v7#:~:text=styles%20graphiques\)\)](https://www.blogdumoderateur.com/test-midjourney-v7#:~:text=styles%20graphiques)))1】. Cette nouvelle version se distingue par un rendu encore amélioré en **photoréalisme** et dans l'imitation de styles artistiques variés ([Test de Midjourney v7 : un modèle bourré de qualités... mais aussi de défauts](https://www.blogdumoderateur.com/test-midjourney-v7#:~:text=styles%20graphiques))) ([https://www.blogdumoderateur.com/test-midjourney-v7#:~:text=styles%20graphiques\)\)](https://www.blogdumoderateur.com/test-midjourney-v7#:~:text=styles%20graphiques)))0】. Les premiers tests soulignent la capacité du modèle à produire des images d'une qualité quasi photographique, gérant bien mieux qu'avant les détails fins (textures, visages) et la cohérence de style. Midjourney 7 innove aussi dans son interaction utilisateur : il introduit un système de **personnalisation** où l'utilisateur peut "former" un profil esthétique. Concrètement, l'abonné doit *liker* ou noter un ensemble d'images proposées, afin que l'IA cerne ses préférences (par exemple, tel style de dessin, telle palette de couleur ([Test de Midjourney v7 : un modèle bourré de qualités... mais aussi de défauts](https://www.blogdumoderateur.com/test-midjourney-v7#:~:text=Cependant%2C%20il%20existe%20pour%20le.styles%20cr%C3%A9%C3%A9s%20par%20d%E2%80%99autres%20utilisatrices))) ([https://www.blogdumoderateur.com/test-midjourney-v7#:~:text=Cependant%2C%20il%20existe%20pour%20le.styles%20cr%C3%A9%C3%A9s%20par%20d%E2%80%99autres%20utilisatrices\)\)](https://www.blogdumoderateur.com/test-midjourney-v7#:~:text=Cependant%2C%20il%20existe%20pour%20le.styles%20cr%C3%A9%C3%A9s%20par%20d%E2%80%99autres%20utilisatrices)))5】. Au bout de ~200 retours, Midjourney génère un profil unique qui conditionne ensuite les résultats – une forme de *fine-tuning* par feedback simplifié pour l'utilisateur. Ce procédé, bien que fastidieux, permet d'obtenir des créations plus alignées sur les goûts individuels. En revanche, Midjourney reste un service fermé (modèle propriétaire, accès payant uniquement, pas d'API publique annoncée). La V7, très attendue par les artistes et designers, consolide la place de Midjourney comme outil de référence pour la création visuelle assistée par IA. Ses limites (toujours des difficultés sur certaines mains, ou des aberrations sur des requêtes complexes) rappellent que le modèle reste expérimental, mais la communauté salue un « *saut qualitatif impressionnant* » dans la continuité de l'évolution fulgurante de Midjourney depuis 20 ([Test de Midjourney v7 : un modèle bourré de qualités... mais aussi de défauts](https://www.blogdumoderateur.com/test-midjourney-v7#:~:text=Plus%20d%E2%80%99un%20an%20apr%C3%A8s%20la.avec%20encore%20de%20nombreux%20d%C3%A9fauts))) ([https://www.blogdumoderateur.com/test-midjourney-v7#:~:text=Plus%20d%E2%80%99un%20an%20apr%C3%A8s%20la.avec%20encore%20de%20nombreux%20d%C3%A9fauts\)\)](https://www.blogdumoderateur.com/test-midjourney-v7#:~:text=Plus%20d%E2%80%99un%20an%20apr%C3%A8s%20la.avec%20encore%20de%20nombreux%20d%C3%A9fauts)))5】).
- **Stable Diffusion – innovations de Stability AI** : Stability AI, pionnier de l'open source visuel, a élargi le champ de la génération au-delà de l'image fixe. Le 18 mars, l'entreprise a annoncé **Stable Virtual Camera**, un nouveau modèle génératif capable de transformer une ou plusieurs images 2D en une **vidéo 3D immersive** ([Stability AI's new AI model turns photos into 3D scenes](https://techcrunch.com/2025/03/18/stability-ai-new-ai-model-turns-photos-into-3d-scenes/#:~:text=Stability%20AI%20has%20released%20a.with%20realistic%20depth%20and%20perspective))) | TechCrunch ([https://techcrunch.com/2025/03/18/stability-ai-new-ai-model-turns-photos-into-3d-scenes/#:~:text=Stability%20AI%20has%20released%20a.with%20realistic%20depth%20and%20perspective\)\)](https://techcrunch.com/2025/03/18/stability-ai-new-ai-model-turns-photos-into-3d-scenes/#:~:text=Stability%20AI%20has%20released%20a.with%20realistic%20depth%20and%20perspective)))6】. Concrètement, à partir d'une simple photo, le modèle génère des *vues nouvelles* de la scène avec des changements d'angle et de perspective réalistes, comme si l'on déplaçait une caméra virtuelle autour de la scène ([Stability AI's new AI model turns photos into 3D scenes](https://techcrunch.com/2025/03/18/stability-ai-new-ai-model-turns-photos-into-3d-scenes/#:~:text=Stable%20Virtual%20Camera%20generates%20%E2%80%99Cnovel,%E2%80%99D))) | TechCrunch ([https://techcrunch.com/2025/03/18/stability-ai-new-ai-model-turns-photos-into-3d-scenes/#:~:text=Stable%20Virtual%20Camera%20generates%20%E2%80%99Cnovel,%E2%80%99D\)\)](https://techcrunch.com/2025/03/18/stability-ai-new-ai-model-turns-photos-into-3d-scenes/#:~:text=Stable%20Virtual%20Camera%20generates%20%E2%80%99Cnovel,%E2%80%99D)))3】. L'utilisateur peut définir une trajectoire de caméra (panoramique, travelling "dolly zoom", mouvement spiralé, etc.) et Stable Virtual Camera produit une séquence vidéo fluide correspondant à ce parcours ([Stability AI's new AI model turns photos into 3D scenes](https://techcrunch.com/2025/03/18/stability-ai-new-ai-model-turns-photos-into-3d-scenes/#:~:text=Stable%20Virtual%20Camera%20generates%20%E2%80%99Cnovel,%E2%80%99D))) | TechCrunch ([https://techcrunch.com/2025/03/18/stability-ai-new-ai-model-turns-photos-into-3d-scenes/#:~:text=Stable%20Virtual%20Camera%20generates%20%E2%80%99Cnovel,%E2%80%99D\)\)](https://techcrunch.com/2025/03/18/stability-ai-new-ai-model-turns-photos-into-3d-scenes/#:~:text=Stable%20Virtual%20Camera%20generates%20%E2%80%99Cnovel,%E2%80%99D)))3】. La version actuelle, en aperçu recherche, permet de créer jusqu'à **1000 frames** en format carré, portrait ou paysa ([Stability AI's new AI model turns photos into 3D scenes](https://techcrunch.com/2025/03/18/stability-ai-new-ai-model-turns-photos-into-3d-scenes/#:~:text=The%20current%20version%20of%20Stable.or%20%E2%80%99Cdynamic%20textures%E2%80%99D%20like%20water))) | TechCrunch ([https://techcrunch.com/2025/03/18/stability-ai-new-ai-model-turns-photos-into-3d-scenes/#:~:text=The%20current%20version%20of%20Stable.or%20%E2%80%99Cdynamic%20textures%E2%80%99D%20like%20water\)\)](https://techcrunch.com/2025/03/18/stability-ai-new-ai-model-turns-photos-into-3d-scenes/#:~:text=The%20current%20version%20of%20Stable.or%20%E2%80%99Cdynamic%20textures%E2%80%99D%20like%20water)))0】. C'est une avancée notable vers la **vidéosynthèse** 3D (sans avoir à modéliser manuellement la scène). Les résultats, bien que prometteurs, sont variables : le modèle peut introduire des artefacts de **flickering** (scintillement) ou des incohérences si la scène est trop ambiguë ou comporte des éléments mouvants comme de l'eau ([Stability AI's new AI model turns photos into 3D scenes](https://techcrunch.com/2025/03/18/stability-ai-new-ai-model-turns-photos-into-3d-scenes/#:~:text=to%201%20C000%20frames%20in%20length.or%20%E2%80%99Cdynamic%20textures%E2%80%99D%20like%20water))) | TechCrunch ([https://techcrunch.com/2025/03/18/stability-ai-new-ai-model-turns-photos-into-3d-scenes/#:~:text=to%201%20C000%20frames%20in%20length.or%20%E2%80%99Cdynamic%20textures%E2%80%99D%20like%20water\)\)](https://techcrunch.com/2025/03/18/stability-ai-new-ai-model-turns-photos-into-3d-scenes/#:~:text=to%201%20C000%20frames%20in%20length.or%20%E2%80%99Cdynamic%20textures%E2%80%99D%20like%20water)))3】. Les visages humains et animaux restent aussi délicats. Stability AI a publié le modèle en **open source (licence non commerciale)** sur Hugging Face afin que la communauté l'expérimente et l'améliore ([Stability AI's new AI model turns photos into 3D scenes](https://techcrunch.com/2025/03/18/stability-ai-new-ai-model-turns-photos-into-3d-scenes/#:~:text=%E2%80%99CHighly%20ambiguous%20scenes%2C%20complex%20camera,%E2%80%99D))) | TechCrunch ([https://techcrunch.com/2025/03/18/stability-ai-new-ai-model-turns-photos-into-3d-scenes/#:~:text=%E2%80%99CHighly%20ambiguous%20scenes%2C%20complex%20camera,%E2%80%99D\)\)](https://techcrunch.com/2025/03/18/stability-ai-new-ai-model-turns-photos-into-3d-scenes/#:~:text=%E2%80%99CHighly%20ambiguous%20scenes%2C%20complex%20camera,%E2%80%99D)))6】. Cette initiative s'inscrit dans la stratégie de Stability AI de se diversifier au-delà de Stable Diffusion : l'entreprise développe parallèlement des modèles pour l'audio (génération de musique), le texte (StableLM) et désormais la vidéo/3D. En coulisses, Stability AI a connu des bouleversements (départ de son CEO fondateur fin 2024, refinancements), et mise sur des innovations comme Stable Virtual Camera pour se démarquer face aux grands concurrents. L'accueil est positif, mais la question de la rentabilité de ces projets open source reste posée, dans un contexte où l'accès gratuit rencontre des coûts d'infrastructure élevés.
- **Runway – Génération vidéo Gen-4** : La startup new-yorkaise Runway, spécialisée dans les outils créatifs dopés à l'IA, a dévoilé le 31 mars son nouveau modèle **Runway Gen-4** ([Runway releases an impressive new video-generating AI model](https://techcrunch.com/2025/03/31/runway-releases-an-impressive-new-video-generating-ai-model/#:~:text=Kyle%20Wiggers))) | TechCrunch ([https://techcrunch.com/2025/03/31/runway-releases-an-impressive-new-video-generating-ai-model/#:~:text=Kyle%20Wiggers\)\)](https://techcrunch.com/2025/03/31/runway-releases-an-impressive-new-video-generating-ai-model/#:~:text=Kyle%20Wiggers)))9】. Succédant à Gen-2 (vidéo à partir de texte) et Gen-3 (version alpha intermédiaire), Gen-4 est présenté comme « l'un des générateurs vidéo IA les plus fidèles qui soient » ([Runway releases an impressive new video-generating AI model](https://techcrunch.com/2025/03/31/runway-releases-an-impressive-new-video-generating-ai-model/#:~:text=AI%20startup%20Runway%20on%20Monday,powered%20video%20generators%20yet))) | TechCrunch ([https://techcrunch.com/2025/03/31/runway-releases-an-impressive-new-video-generating-ai-model/#:~:text=AI%20startup%20Runway%20on%20Monday,powered%20video%20generators%20yet\)\)](https://techcrunch.com/2025/03/31/runway-releases-an-impressive-new-video-generating-ai-model/#:~:text=AI%20startup%20Runway%20on%20Monday,powered%20video%20generators%20yet)))6】. Les avancées mises en avant : une **cohérence améliorée** sur la durée de la vidéo, avec maintien du même personnage, décor et objets à travers des plans successifs ([Runway releases an impressive new video-generating AI model](https://techcrunch.com/2025/03/31/runway-releases-an-impressive-new-video-generating-ai-model/#:~:text=Called%20Gen,perspectives%20and%20positions%20within%20scenes))) | TechCrunch ([https://techcrunch.com/2025/03/31/runway-releases-an-impressive-new-video-generating-ai-model/#:~:text=Called%20Gen,perspectives%20and%20positions%20within%20scenes\)\)](https://techcrunch.com/2025/03/31/runway-releases-an-impressive-new-video-generating-ai-model/#:~:text=Called%20Gen,perspectives%20and%20positions%20within%20scenes)))8】. Là où les modèles précédents avaient du mal à garder un visage ou un personnage identique d'une scène à l'autre, Gen-4 y parvient nettement mieux. De même, il gère des *"world environments"* consistants : on peut demander une séquence dans un univers particulier, et l'IA conservera le style et les éléments de ce monde tout du long ([Runway releases an impressive new video-generating AI model](https://techcrunch.com/2025/03/31/runway-releases-an-impressive-new-video-generating-ai-model/#:~:text=Called%20Gen,perspectives%20and%20positions%20within%20scenes))) | TechCrunch ([https://techcrunch.com/2025/03/31/runway-releases-an-impressive-new-video-generating-ai-model/#:~:text=Called%20Gen,perspectives%20and%20positions%20within%20scenes\)\)](https://techcrunch.com/2025/03/31/runway-releases-an-impressive-new-video-generating-ai-model/#:~:text=Called%20Gen,perspectives%20and%20positions%20within%20scenes)))9】. Runway indique que Gen-4 sait utiliser des **références visuelles** fournies par l'utilisateur : par exemple on peut donner la photo d'un acteur réel, et la vidéo générée mettra en scène un personnage ressemblant à cette référence dans différentes poses et éclairages ([Runway releases an impressive new video-generating AI model](https://techcrunch.com/2025/03/31/runway-releases-an-impressive-new-video-generating-ai-model/#:~:text=Runway%2C%20which%20is%20C2%A0backed%20C2%A0by%20investors.generated%20video))) | TechCrunch ([https://techcrunch.com/2025/03/31/runway-releases-an-impressive-new-video-generating-ai-model/#:~:text=Runway%2C%20which%20is%20C2%A0backed%20C2%A0by%20investors.generated%20video\)\)](https://techcrunch.com/2025/03/31/runway-releases-an-impressive-new-video-generating-ai-model/#:~:text=Runway%2C%20which%20is%20C2%A0backed%20C2%A0by%20investors.generated%20video)))7】. Le tout **sans fine-tuning** nécessaire, via de simples prompts ou uploads d'image ([Runway releases an impressive new video-generating AI model](https://techcrunch.com/2025/03/31/runway-releases-an-impressive-new-video-generating-ai-model/#:~:text=Runway%2C%20which%20is%20C2%A0backed%20C2%A0by%20investors.generated%20video))) | TechCrunch ([https://techcrunch.com/2025/03/31/runway-releases-an-impressive-new-video-generating-ai-model/#:~:text=Runway%2C%20which%20is%20C2%A0backed%20C2%A0by%20investors.generated%20video\)\)](https://techcrunch.com/2025/03/31/runway-releases-an-impressive-new-video-generating-ai-model/#:~:text=Runway%2C%20which%20is%20C2%A0backed%20C2%A0by%20investors.generated%20video)))7】).

• **Autres nouveautés visuelles** : On note également en mars des efforts sur la **génération multimodale** plus large. Par exemple, OpenAI a enrichi son ChatGPT d'une fonction de **génération d'images** intégrée (basée sur DALL·E3) pour les utilisateurs Plus. Cette fonctionnalité a connu un tel engouement début mars que les serveurs ont peiné à suivre : Sam Altman a indiqué « *nos GPU sont en train de fondre* » en parlant de l'explosion des requêtes images, forçant OpenAI à limiter provisoirement le nombre d'images générées par utilisateur ([OpenAI says 'our GPUs are melting' as it limits ChatGPT image generation requests | The Verge](https://www.theverge.com/news/637542/chatgpt-says-our-gpus-are-melting-as-it-puts-limit-on-image-generation-requests#%7E:text=The%20error%20around%20ChatGPT%E2%80%99s%20more,handling%20the%20avalanche%20of%20requests))) ([https://www.theverge.com/news/637542/chatgpt-says-our-gpus-are-melting-as-it-puts-limit-on-image-generation-requests#%7E:text=The%20demand%20crunch%20already%20caused,to%20three%20images%20per%20day\)\)5](https://www.theverge.com/news/637542/chatgpt-says-our-gpus-are-melting-as-it-puts-limit-on-image-generation-requests#%7E:text=The%20demand%20crunch%20already%20caused,to%20three%20images%20per%20day))5)). Cela montre l'appétit du public pour les IA créatives (ici, des millions d'images de style *Studio Ghibli* générées en quelques jou ([OpenAI says 'our GPUs are melting' as it limits ChatGPT image generation requests | The Verge](https://www.theverge.com/news/637542/chatgpt-says-our-gpus-are-melting-as-it-puts-limit-on-image-generation-requests#%7E:text=The%20demand%20crunch%20already%20caused,to%20three%20images%20per%20day))5) ([https://www.theverge.com/news/637542/chatgpt-says-our-gpus-are-melting-as-it-puts-limit-on-image-generation-requests#%7E:text=The%20demand%20crunch%20already%20caused,to%20three%20images%20per%20day\)\)5](https://www.theverge.com/news/637542/chatgpt-says-our-gpus-are-melting-as-it-puts-limit-on-image-generation-requests#%7E:text=The%20demand%20crunch%20already%20caused,to%20three%20images%20per%20day))5))) et rappelle le **coût matériel** colossal de ces modèles. De même, sur le front de la **3D**, NVIDIA a annoncé à la GTC de nouveaux modèles *Cosmos* capables de générer des **mondes virtuels synthétiques** pour entraîner des IA robotique ([GTC 2025 – Announcements and Live Updates | NVIDIA Blog](https://blogs.nvidia.com/blog/nvidia-keynote-at-gtc-2025-ai-news-live-updates/?utm_source=chatgpt.com#%7E:text=NVIDIA%20also%20announced%20a%20major,unprecedented%20control%20over%20world%20generation))4) ([https://blogs.nvidia.com/blog/nvidia-keynote-at-gtc-2025-ai-news-live-updates/?utm_source=chatgpt.com#%7E:text=NVIDIA%20also%20announced%20a%20major,unprecedented%20control%20over%20world%20generation\)\)4](https://blogs.nvidia.com/blog/nvidia-keynote-at-gtc-2025-ai-news-live-updates/?utm_source=chatgpt.com#%7E:text=NVIDIA%20also%20announced%20a%20major,unprecedented%20control%20over%20world%20generation))4)). L'idée est de créer, via l'IA, un nombre infini d'environnements de simulation dans lesquels les robots pourront apprendre, plutôt que de se limiter aux données du monde réel ([GTC 2025 – Announcements and Live Updates | NVIDIA Blog](https://blogs.nvidia.com/blog/nvidia-keynote-at-gtc-2025-ai-news-live-updates/?utm_source=chatgpt.com#%7E:text=E%2%80%9CUsing%20Omniverse%20to%20condition%20Cosmos%2C,the%20same%20time%2C%E2%80%9D%20Huang%20said))) ([https://blogs.nvidia.com/blog/nvidia-keynote-at-gtc-2025-ai-news-live-updates/?utm_source=chatgpt.com#%7E:text=E%2%80%9CUsing%20Omniverse%20to%20condition%20Cosmos%2C,the%20same%20time%2C%E2%80%9D%20Huang%20said\)\)](https://blogs.nvidia.com/blog/nvidia-keynote-at-gtc-2025-ai-news-live-updates/?utm_source=chatgpt.com#%7E:text=E%2%80%9CUsing%20Omniverse%20to%20condition%20Cosmos%2C,the%20same%20time%2C%E2%80%9D%20Huang%20said)))); Ces progrès en vision, vidéo et simulation confirment que l'IA générative ne se limite plus au texte : elle envahit la création sous toutes ses formes (image fixe, art, vidéo, 3D interactive), ouvrant des perspectives dans le divertissement, l'éducation, la conception produit, etc. tout en posant de nouveaux défis (véracité des contenus, propriété intellectuelle, empreinte carbone du calcul intensif requis).

- **Évolutions des frameworks ML** : Les bibliothèques et frameworks majeurs continuent de s'adapter à l'ère des LLM. **PyTorch** (meta/OSS) et **TensorFlow** (Google) ont tous deux publié des mises à jour focalisées sur la **performance à grande échelle** et l'optimisation des déploiements de modèles géants. PyTorch 2.1 (en cours) introduit par exemple des améliorations du compilateur *torch.compile* pour accélérer l'inférence des LLM sur GPU, et intègre mieux les optimisations CUDA les plus récentes d'après NVI ([The Performance of CUDA with the Flexibility of PyTorch - NVIDIA \(https://www.nvidia.com/en-us/on-demand/session/gtc25-S71946/#%7E:text=NVIDIA%20www.Date%3A%20March%202025\),L9](https://www.nvidia.com/en-us/on-demand/session/gtc25-S71946/#%7E:text=NVIDIA%20www.Date%3A%20March%202025),L9)). Du côté de TensorFlow, l'accent est mis sur la **portabilité** (exécutions optimisées sur TPU v5, CPU, GPU via XLA) et sur le support de longs contextes (ops spéciales pour gérer des séquences de >100k tokens sans exploser la mémoire). **JAX**, la librairie de calcul automatique de Google, reste très utilisée en recherche (par ex. le projet *MaxText* utilise JAX sur TPU v5 pour entraîner des LLM de manière *Auto-scala* ([JAX things to watch for in 2025 - by Grigory Sapunov - Gonzo ML \(https://gonzoml.substack.com/p/jax-things-to-watch-for-in-2025#%7E:text=ML%20gonzoml.GPUs%20for%20training%20and%20inference\),L8](https://arxiv.org/abs/2502.11804))). On observe aussi un rapprochement des écosystèmes : PyTorch intègre des contributions Nvidia (framework NeMo), tandis que TensorFlow coopère avec JAX (interopérabilité via TF-JAX). Si aucune version « 3.0 » de ces frameworks n'est sortie en mars, leurs *roadmaps* indiquent une évolution continue pour supporter l'IA générative en production : amélioration de la **distribution des modèles sur plusieurs GPUs**, réduction des coûts mémoire (quantification plus agressive), et insertion de garde-fous éthiques (certains proposent désormais des plugins pour filtrer les outputs offensants en standard).
- **LangChain et agents LLM** : LangChain, le framework open source de référence pour construire des applications autour des LLM, a connu plusieurs mises à jour significatives en mars 2025. Constatant l'essor des **agents IA** (systèmes où plusieurs modèles collaborent pour accomplir une tâche), l'équipe a lancé **LangGraph**, une bibliothèque dédiée à la création d'agents multi-LM sophistiqués. La version *LangGraph 0.3* apporte notamment des *agents préconstruits* pour aider les développeurs à démarrer plus v ([LangChain - Changelog \(https://changelog.langchain.com/?page=2#%7E:text=Key%20features%20include%3A%20The%20library,26\),64](https://changelog.langchain.com/?page=2#%7E:text=Key%20features%20include%3A%20The%20library,26),64)). Surtout, le module **LangGraph Swarm** permet d'orchestrer un essaim d'agents IA en parallèle – par exemple une dizaine de petits modèles se répartissant un problème en sous-tâches et communiquant entre eux pour converger vers une solut ([LangChain - Changelog \(https://changelog.langchain.com/?page=2#%7E:text=March%202025\),62](https://changelog.langchain.com/?page=2#%7E:text=March%202025),62)). En outre, LangChain a ajouté le support natif de **Claude 3.7 d'Anthropic** (via SDK Python/JS) afin de profiter des avancées de ce mod ([LangChain - Changelog \(https://changelog.langchain.com/?page=2#%7E:text=anthropic%20giving%20full%20control%20over.March%201%20C%202025\),64](https://changelog.langchain.com/?page=2#%7E:text=anthropic%20giving%20full%20control%20over.March%201%20C%202025),64)). Un autre ajout marquant est le *LangMem SDK*, bibliothèque introduite mi-mars pour doter les agents d'une **mémoire long-terme** persista ([LangMem SDK for agent long-term memory - LangChain Blog \(https://blog.langchain.dev/langmem-sdk-launch/#%7E:text=LangMem%20SDK%20for%20agent%20long-term-memory%20memory\),13](https://blog.langchain.dev/langmem-sdk-launch/#%7E:text=LangMem%20SDK%20for%20agent%20long-term-memory%20memory),13)). Cela permet par exemple à un agent conversationnel de retenir des faits sur l'utilisateur ou des événements passés et d'adapter son comportement sur la durée (apprentissage en continu). Cette course aux *agents intelligents* reflète le buzz autour des AutoGPT et autres assistants autonomes : LangChain vise à fournir l'infrastructure outillée pour en développer de manière robuste, sans repartir de zéro. Néanmoins, la communauté pointe que l'empilement de tels agents complexifie fortement le débogage (« *rabbit hole* » selon certains développeur ([langchain is still a rabbit hole in 2025 : r/LocalLLaMA - Reddit \(https://www.reddit.com/r/LocalLLaMA/comments/1iuda08/langchain_is_still_a_rabbit_hole_in_2025/#%7E:text=langchain%20is%20still%20a%20rabbit,is%20the%20case%20that%20we%20are%20using%20a%20lot%20of%20tools%20to%20debug%20these%20agents,1](https://www.reddit.com/r/LocalLLaMA/comments/1iuda08/langchain_is_still_a_rabbit_hole_in_2025/#%7E:text=langchain%20is%20still%20a%20rabbit,is%20the%20case%20that%20we%20are%20using%20a%20lot%20of%20tools%20to%20debug%20these%20agents,1))).
- **Hugging Face Hub & co** : La plateforme Hugging Face demeure au centre de l'écosystème open source, avec plus de 250 000 modèles accessibles. En mars, HF a annoncé plusieurs partenariats stratégiques pour renforcer la **fiabilité** et la **sécurité** de cet écosystème foisonnant. Lors des *MLOps Days* de New York (début mars), la

société de sécurité JFrog a révélé une alliance avec Hugging Face visant à analyser **automatiquement les modèles du Hub à la recherche de failles ou de comportements malveillant (JFrog s'associe à Hugging Face pour assurer la sécurité des modèles GenAI (<https://www.solutions-numeriques.com/jfrog-sassocie-a-hugging-face-pour-assurer-la-securite-des-modeles-genai/#%7E:text=Le%20sp%C3%A9cialiste%20de%20la%20s%C3%A9curit%C3%A9,vuln%C3%A9rabilit%C3%A9s%20et%20de%20mod%C3%A8s%20malveillants>)) (JFrog s'associe à Hugging Face pour assurer la sécurité des modèles GenAI (<https://www.solutions-numeriques.com/jfrog-sassocie-a-hugging-face-pour-assurer-la-securite-des-modeles-genai/#%7E:text=Cette%20int%C3%A9gration%20consistera%2C%20sch%C3%A9matiquement%2C%20%C3%A0%20Karas%2C%20CTO%20de%20JFrog%20Security>))33】.

Concrètement, chaque nouveau modèle uploadé pourra être scanné pour détecter d'éventuels logiciels espions intégrés, fuites de données sensibles ou biais indésirables. Cette initiative répond à une préoccupation croissante : à mesure que les modèles ML sont utilisés en entreprise, garantir leur **intégrité** devient aussi important que la sécurité du code logiciel. Par ailleurs, Hugging Face continue d'élargir son champ d'action. Après le texte et l'image, la société s'attaque à la **robotique open source** : elle a officiellement acquis en mars la startup française *Pollen Robotics*, connue pour son robot humanoïde modulaire *Reac (Hugging Face acquiert Pollen Robotics (<https://www.planeterobots.com/2025/04/14/hugging-face-acquiert-pollen-robotics/#%7E:text=Hugging%20Face%20rach%C3%A8te%20Pollen%20Robotics>))66】. En s'appuyant sur la bibliothèque Robotique *LeRobot* (lancée en 2024) et sur le hub de modèles, HF ambitionne de bâtir la plateforme ouverte de référence pour la **robotique intelligente** – un écosystème où matériels open source (bras robotique low-cost SO-100, etc.) et modèles IA entraînés (comme *GR00T N1* de NVIDIA, voir plus bas) se rencontrent (Hugging Face acquiert Pollen Robotics (<https://www.planeterobots.com/2025/04/14/hugging-face-acquiert-pollen-robotics/#%7E:text=accessible%20du%20march%C3%A9%2C%20et%20l'E2%80%99un,et%20la%20v%C3%A9rification%20de%20LeRobot>))07】. L'acquisition de Pollen s'inscrit dans cette vision : Thomas Wolf (HF) évoque « *la robotique comme prochaine frontière de l'IA, qui doit être ouverte et abordable* », invitant la communauté à construire des assistants physiques aussi facilement que des chatbots (Hugging Face acquiert Pollen Robotics (<https://www.planeterobots.com/2025/04/14/hugging-face-acquiert-pollen-robotics/#%7E:text=Thomas%20Wolf%2C%20cofondateur%20et%20directeur%20Port%C3%A9s%20par%20la>)) (Hugging Face acquiert Pollen Robotics (<https://www.planeterobots.com/2025/04/14/hugging-face-acquiert-pollen-robotics/#%7E:text=Nous%20pensons%20que%20la%20robotique,Face%20est%20un%20lieu%20naturel>))79】. Enfin, sur le volet **collaborations cloud**, on note qu'IBM a intégré les nouveaux modèles Llama 3.1 et Mistral 3.1 sur sa plateforme Watsonx, et que NVIDIA collabore avec HF (depuis nov. 2024) pour accélérer l'entraînement des modèles du Hub sur ses (Hugging Face acquiert Pollen Robotics (<https://www.planeterobots.com/2025/04/14/hugging-face-acquiert-pollen-robotics/#%7E:text=,et%20la%20v%C3%A9rification%20de%20LeRobot>))04】. Hugging Face s'affirme ainsi comme un carrefour incontournable réunissant modèles, données et désormais matériel, tout en s'assurant que la confiance (via la sécurité, l'éthique) soit au rendez-vous.

- **Outils de déploiement et de gouvernance** : En marge des frameworks, plusieurs outils et pratiques émergent pour accompagner le déploiement industriel de l'IA. Par exemple, **Ray** et **DeepSpeed** (outils d'orchestration et d'optimisation distribuée) ont publié des guides en mars pour faciliter l'hébergement de LLM avec contraintes de coût réduites. **Scikit-learn** (bibliothèque ML classique) a sorti sa v1.6.1 en janvier et prépare scikit-learn 1.7 pour 2025, incluant quelques algorithmes de ML plus efficaces et une compatibilité améliorée avec le GPU via C ([scikit-learn: machine learning in Python — scikit-learn 1.6.1 ...](https://scikit-learn.org/#%7E:text=learn.org%20On%20learn%201.6) (<https://scikit-learn.org/#%7E:text=learn.org%20On%20learn%201.6>))L8】. Du côté **MLOps**, outre la sécurité (ex : partenariat HF-JFrog cité), on discute beaucoup de la surveillance continue des modèles en production. Des conférences comme les MLOps Days et des livres blancs parus ce mois-ci soulignent les bonnes pratiques pour détecter les **dérives de performance** d'un modèle après déploiement (concept *Data Drift*), ou pour auditer les biais de manière régulière. On a aussi vu apparaître des solutions commerciales de *"Model as a Service"* locales : certaines startups proposent aux entreprises d'héberger sur site leur propre LLM finetuné, pour combiner confidentialité et performance, un compromis qui gagne en popularité face aux seules API cloud. Enfin, la **standardisation** commence : l'ISO et le IEEE travaillent sur des normes de documentation des modèles IA (cartes de score, méta-données sur les données d'entraînement) afin de faciliter l'échange d'informations entre développeurs, régulateurs et utilisateurs. En somme, l'écosystème outillage s'adapte à la fois *techniquement* (pour supporter des modèles toujours plus gros, multimodaux, en inférence temps réel) et *organisationnellement* (pour intégrer l'IA de manière sécurisée et responsable dans les pipelines existants).

□ Faits marquants recherche, applications et éthique

- **Manus AI : agent autonome** – Une démonstration marquante est venue de Chine le 6 mars : la startup Monica AI a dévoilé **Manus AI**, présenté comme un agent conversationnel **totalement autonome (IA : 3 grandes avancées en mars 2025 qui vont tout changer - Upskilling (<https://upskilling.com/ia-3-grandes-avancees-en-mars-2025-qui-vont-tout-changer/#%7E:text=Le%206%20mars%202025%2C%20la%20besoin%20d'E2%80%99une%20intervention%20humaine%20constante>))L45】. Contrairement aux chatbots traditionnels qui répondent simplement aux requêtes, Manus est conçu pour **prendre des décisions par lui-même** et exécuter des tâches complexes de bout en bout, avec un minimum d'intervention humaine. Par exemple, il peut analyser une situation client, déterminer les actions à mener (envoyer un remboursement, planifier un appel, etc.) et les réaliser automatiquement (IA : 3 grandes avancées en mars 2025 qui vont tout changer - Upskilling (<https://upskilling.com/ia-3-grandes-avancees-en-mars-2025-qui-vont-tout-changer/#%7E:text=chatbots%20classiques%2C%20A0Manus%2C%20est%20capable%20de%20fonctionner, besoin%20d'E2%80%99une%20intervention%20humaine%20constante>))L46】. Plusieurs plateformes e-commerce chinoises expérimentent déjà Manus pour gérer le support client de manière automatisée : il pourrait « *fluidifier les échanges et améliorer l'expérience utilisateur* » selon ses créateurs (IA : 3 grandes avancées en mars 2025 qui vont tout changer - Upskilling (<https://upskilling.com/ia-3-grandes-avancees-en-mars-2025-qui-vont-tout-changer/#%7E:text=Son%20potentiel%20est%20immense,am%C3%A9liorer%20l'E2%80%99exp%C3%A9rience%20utilisateur%20et%20fluidifier>))L47】. Le potentiel en entreprise est énorme (automatisation de processus administratifs, optimisation logistique...), mais ce prototype soulève aussi des questions cruciales. Jusqu'où peut-on déléguer des décisions à une IA sans supervision ? Manus relance le débat sur l'**autonomie des IA** : bénéfices (gain d'efficacité) contre risques (erreurs non détectées, manque de contrôle humain (Intelligence Artificielle : Les 5 actualités majeures du 14 mars 2025 (<https://fr.linkedin.com/pulse/intelligence-artificielle-les-5-actualit%C3%A9s-majeures-du-nobili-mpeve/#%7E:text=La%20Chine%20a%20d%C3%A9voilé%20C3%A9%20Manus,pour%20encadrer%20ces%20nouvelles%20capacit%C3%A9s>)) (Intelligence Artificielle : Les 5 actualités majeures du 14 mars 2025 (<https://fr.linkedin.com/pulse/intelligence-artificielle-les-5-actualit%C3%A9s-majeures-du-nobili-mpeve/#%7E:text=inqui%C3%A9tudes%203A%20jusqu'E2%80%99o%C3%B9%20peut%20aller,pour%20encadrer%20ces%20nouvelles%20capacit%C3%A9s>))L94】. Les observateurs notent qu'une **régulation stricte** sera sans doute nécessaire avant d'envisager un déploiement à grande échelle de tels agents autonomes. En tout cas, cette annonce illustre que sur le plan technologique, la barrière de l'autonomie est en train d'être franchie expérimentalement, ce qui était un objectif de longue date en IA (*agents intelligents capables d'agir dans un environnement*).
- **Sommet mondial sur l'IA à Paris** – Les 11–12 mars, la France a accueilli à Paris le *Sommet de l'Action pour l'Intelligence Artificielle*, réunissant des représentants de **60 États** et des centaines d'experts (IA : 3 grandes avancées en mars 2025 qui vont tout changer - Upskilling (<https://upskilling.com/ia-3-grandes-avancees-en-mars-2025-qui-vont-tout-changer/#%7E:text=Un%20sommet%20mondial%20pour%20encadrer,l'E2%80%99intelligence%20artificielle>))L68】. L'objectif était de discuter d'un cadre international pour une IA plus **éthique, inclusive et durable**. Durant ce sommet, une **déclaration commune** a été signée, posant des principes de transparence des algorithmes, de respect de la vie privée et d'**absence de biais** dans les systèmes (IA : 3 grandes avancées en mars 2025 qui vont tout changer - Upskilling

([https://upskilling.com/ia-3-grandes-avancees-en-mars-2025-qui-vont-tout-](https://upskilling.com/ia-3-grandes-avancees-en-mars-2025-qui-vont-tout-changer/#%7E:text=La%20mont%C3%A9e%20en%20puissance%20des,IA%20plus%20transparente%20et%20responsable)

[change/#%7E:text=La%20mont%C3%A9e%20en%20puissance%20des,IA%20plus%20transparente%20et%20responsable](https://upskilling.com/ia-3-grandes-avancees-en-mars-2025-qui-vont-tout-changer/#%7E:text=La%20mont%C3%A9e%20en%20puissance%20des,IA%20plus%20transparente%20et%20responsable))L75】. Toutefois, les approches divergent : l'Union Européenne a plaidé pour un encadrement juridique strict et rapide (dans la lignée de l'AI Act européen qui pourrait entrer en vigueur fin 2025), tandis que les États-Unis et la Chine se montrent plus flexibles, privilégiant l'innovation et l'autorégulation sur certains asp ([IA : 3 grandes avancées en mars 2025 qui vont tout changer - Upskilling](https://upskilling.com/ia-3-grandes-avancees-en-mars-2025-qui-vont-tout-changer/#%7E:text=ou%20l%E2%80%99automatisation%20excessive%20des%20emplois,IA%20plus%20transparente%20et%20responsable) ([https://upskilling.com/ia-3-grandes-avancees-en-mars-2025-qui-vont-tout-](https://upskilling.com/ia-3-grandes-avancees-en-mars-2025-qui-vont-tout-changer/#%7E:text=ou%20l%E2%80%99automatisation%20excessive%20des%20emplois,IA%20plus%20transparente%20et%20responsable)

[change/#%7E:text=ou%20l%E2%80%99automatisation%20excessive%20des%20emplois,IA%20plus%20transparente%20et%20responsable](https://upskilling.com/ia-3-grandes-avancees-en-mars-2025-qui-vont-tout-changer/#%7E:text=ou%20l%E2%80%99automatisation%20excessive%20des%20emplois,IA%20plus%20transparente%20et%20responsable))L75】. Ce sommet fait écho à d'autres initiatives (Plan d'action du G7 sur l'IA, forum GPAI, etc.) et traduit l'urgence ressentie par les gouvernements de **se saisir du sujet**. Pour les entreprises, cela annonce un futur proche où il sera indispensable de prouver la conformité de leurs systèmes d'IA (traçabilité des données d'entraînement, audits indépendants...) afin de répondre aux normes à v ([IA : 3 grandes avancées en mars 2025 qui vont tout changer - Upskilling](https://upskilling.com/ia-3-grandes-avancees-en-mars-2025-qui-vont-tout-changer/#%7E:text=Selon%20une%20enqu%C3%AAte%20du%20cabinet,Donc%20essentiel%20pour%20rester%20comp%C3%A9titif) ([https://upskilling.com/ia-3-grandes-avancees-en-mars-2025-qui-vont-tout-](https://upskilling.com/ia-3-grandes-avancees-en-mars-2025-qui-vont-tout-changer/#%7E:text=Selon%20une%20enqu%C3%AAte%20du%20cabinet,Donc%20essentiel%20pour%20rester%20comp%C3%A9titif)

[change/#%7E:text=Selon%20une%20enqu%C3%AAte%20du%20cabinet,Donc%20essentiel%20pour%20rester%20comp%C3%A9titif](https://upskilling.com/ia-3-grandes-avancees-en-mars-2025-qui-vont-tout-changer/#%7E:text=Selon%20une%20enqu%C3%AAte%20du%20cabinet,Donc%20essentiel%20pour%20rester%20comp%C3%A9titif))L87】. Une enquête Deloitte citée indique que 75% des dirigeants anticipent un impact majeur de la régulation IA sur leur business d'ici ([IA : 3 grandes avancées en mars 2025 qui vont tout changer - Upskilling](https://upskilling.com/ia-3-grandes-avancees-en-mars-2025-qui-vont-tout-changer/#%7E:text=Selon%20une%20enqu%C3%AAte%20du%20cabinet,Donc%20essentiel%20pour%20rester%20comp%C3%A9titif) ([https://upskilling.com/ia-3-grandes-avancees-en-mars-2025-qui-vont-tout-](https://upskilling.com/ia-3-grandes-avancees-en-mars-2025-qui-vont-tout-changer/#%7E:text=Selon%20une%20enqu%C3%AAte%20du%20cabinet,Donc%20essentiel%20pour%20rester%20comp%C3%A9titif)

- **Code de bonnes pratiques pour l'IA générative** – En complément de la régulation *hard law*, les acteurs se sont accordés sur des **lignes directrices volontaires**. Début mars, un **code de conduite** européen pour les IA à usage général (fondation models) a été publié, cherchant un équilibre entre contrer les abus et ne pas freiner l'innova ([Intelligence Artificielle : Les 5 actualités majeures du 14 mars 2025](https://fr.linkedin.com/pulse/intelligence-artificielle-les-5-actualit%C3%A9s-majeures-du-nobili-mpeve#%7E:text=%E2%9A%96%E2%8F%20Face%20%C3%A0%20l%E2%80%99essor%20rapide,transparence%2C%20de%20s%C3%A9curit%C3%A9%20et%20d%20Intelligence%20Artificielle%20Les%205%20actualit%C3%A9s%20majeures%20du%2014%20mars%202025) ([https://fr.linkedin.com/pulse/intelligence-artificielle-les-5-actualit%C3%A9s-majeures-du-nobili-](https://fr.linkedin.com/pulse/intelligence-artificielle-les-5-actualit%C3%A9s-majeures-du-nobili-mpeve#%7E:text=%E2%9A%96%E2%8F%20Face%20%C3%A0%20l%E2%80%99essor%20rapide,transparence%2C%20de%20s%C3%A9curit%C3%A9%20et%20d%20Intelligence%20Artificielle%20Les%205%20actualit%C3%A9s%20majeures%20du%2014%20mars%202025)
- **Nvidia GTC 2025 – matériel et robotique** – La conférence annuelle **GTC 2025** de Nvidia (17–21 mars à San Jose) a confirmé le rôle central du matériel dans la révolution ([IA : 3 grandes avancées en mars 2025 qui vont tout changer - Upskilling](https://upskilling.com/ia-3-grandes-avancees-en-mars-2025-qui-vont-tout-changer/#%7E:text=D%C3%A9but%20mars,un%20code%20de%20conduite%20europ%C3%A9en%20pour%20les%20IA%20%C3%A0%20usage%20g%C3%A9n%C3%A9ral%20(fondation%20models)%20a%20%C3%A9t%C3%A9%20publi%C3%A9,cherchant%20un%20%C3%A9quilibre%20entre%20contrer%20les%20abus%20et%20ne%20pas%20freiner%20l'innovation) ([https://upskilling.com/ia-3-grandes-avancees-en-mars-2025-qui-vont-tout-](https://upskilling.com/ia-3-grandes-avancees-en-mars-2025-qui-vont-tout-changer/#%7E:text=D%C3%A9but%20mars,un%20code%20de%20conduite%20europ%C3%A9en%20pour%20les%20IA%20%C3%A0%20usage%20g%C3%A9n%C3%A9ral%20(fondation%20models)%20a%20%C3%A9t%C3%A9%20publi%C3%A9,cherchant%20un%20%C3%A9quilibre%20entre%20contrer%20les%20abus%20et%20ne%20pas%20freiner%20l'innovation)

non seulement comme fournisseur de **GPU** pour l'IA générative, mais aussi comme un acteur clé de la **robotique intelligente**, fournissant à la fois les cerveaux (modèles GR00T/Cosmos), le simulateur (Omniverse/Newton) et le hardware spécialisé (capteurs, Jetson, etc.). Cela augure d'une accélération de la convergence IA + robotique dans l'industrie, la santé (robots chirurgicaux prése ([GTC 2025 – Announcements and Live Updates | NVIDIA Blog \(https://blogs.nvidia.com/blog/nvidia-keynote-at-gtc-2025-ai-news-live-updates/?utm_source=chatgpt.com#:~:text=match%20at%20L895%20Image%3A%20surgical,its%20dexterity%20at%20GTC%202024\)\)L24](#)]), la logistique, etc., avec l'espoir de combler le manque de main d'œuvre à v ([GTC 2025 – Announcements and Live Updates | NVIDIA Blog \(https://blogs.nvidia.com/blog/nvidia-keynote-at-gtc-2025-ai-news-live-updates/?utm_source=chatgpt.com#:~:text=Describing%20robots%20as%20the%20next,generation%20robotics\)\)714](#)] mais en posant aussi la question de la place de **3 milliards de robots humanoïdes d'ici 2060** (chiffre prospectif par Bank of America) dans nos soci ([Intelligence Artificielle : Les 5 actualités majeures du 14 mars 2025 \(https://fr.linkedin.com/pulse/intelligence-artificielle-les-5-actualit%C3%A9s-majeures-du-nobili-mpeve#:~:text=2%E2%8F%B8%E2%83%A3%20L%27IA%20et%20la%20robotique,de%20robots%20humano%C3%AFdes%20d%E2%80%99ici%202060\)\)L70](#)]).

- **IA & Finance** – Un secteur transformé silencieusement par l'IA est la **finance**. En mars, on a constaté que les grandes institutions financières, des banques centrales aux fonds d'investissement, intensifient l'adoption de l'IA pour la **gestion des risques et la prévision économiq ([Intelligence Artificielle : Les 5 actualités majeures du 14 mars 2025 \(https://fr.linkedin.com/pulse/intelligence-artificielle-les-5-actualit%C3%A9s-majeures-du-nobili-mpeve#:~:text=5%E2%8F%B8%E2%83%A3%20L%E2%80%99IA%20dans%20la%20finance,risques%20et%20la%20pr%C3%A9vision%20%C3%A9conomique\)\)104](#)]). Des algorithmes de machine learning avancés analysent des volumes massifs de données de marché, de news, de signaux macroéconomiques afin de détecter précocement les signes de crise ou d'optimiser les allocations d'actifs. Par exemple, certains hedge funds utilisent des modèles génératifs pour simuler des scénarios extrêmes et tester la robustesse de leur portefeuille. Résultat : les prévisions de tendances se font avec une précision accrue et une réactivité supérieure, ce qui peut réduire l'exposition aux chocs financ ([Intelligence Artificielle : Les 5 actualités majeures du 14 mars 2025 \(https://fr.linkedin.com/pulse/intelligence-artificielle-les-5-actualit%C3%A9s-majeures-du-nobili-mpeve#:~:text=Les%20grandes%20institutions%20financi%C3%A8res%20acc%C3%A9l%C3%A8rent,les%20sp%C3%A9culatives%20aliment%C3%A9es%20par%20Toutefois,ce%20financ%20algorithmique,soulève,des,enjeux,de,transparence,les,décisions,prises,par,une,IA,complexe,sont,difficiles,à,expliquer,et,des,craintes,de,**bulles,spéculatives,amplifiées,par,l',\(Intelligence Artificielle : Les 5 actualités majeures du 14 mars 2025 \(https://fr.linkedin.com/pulse/intelligence-artificielle-les-5-actualit%C3%A9s-majeures-du-nobili-mpeve#:~:text=Les%20grandes%20institutions%20financi%C3%A8res%20acc%C3%A9l%C3%A8rent,les%20sp%C3%A9culatives%20aliment%C3%A9es%20par%20Si,de,plus,en,plus,d'agents,autonomes,tradent,sur,les,marchés,pourraient-ils,engendrer,collectivement,des,dynamiques,imprévues,Les,régulateurs,financiers,s'intéressent,de,près,à,ces,questions,cherchant,à,adapter,les,contrôles,de,risque,systémique,à,l'ère,de,l'IA.](#)

- **Recherche académique et sciences** – Dans le monde académique, mars 2025 a apporté son lot de publications remarquables en IA. En **biologie**, des chercheurs français de l'INRAE ont annoncé un **système IA hybride** capable de **concevoir de nouvelles protéines** en combinant apprentissage profond et règles physico-chim ([Major breakthrough in biology: a hybrid generative AI designs new ... \(https://www.inrae.fr/en/news/major-breakthrough-biology-hybrid-generative-ai-designs-new-molecules#:~:text=Major%20breakthrough%20in%20biology%3A%20a,from%20physics%20or%20made\)\)L37](#)]). Cette approche génère des protéines satisfaisant simultanément des critères de stabilité physique et des critères appris sur des bases de données biologiques – une avancée vers des enzymes sur mesure pour l'industrie ou de nouveaux médicaments, où l'IA accélère la découverte. En **physique**, la collaboration Quantinuum a présenté un cadre de *Generative Quantum AI* utilisant des données quantiques aléatoires pour entraîner des modèles hybrides, ouvrant une piste pour résoudre des problèmes d'optimisation complexes via un mélange quantique-class ([Quantinuum Announces Generative Quantum AI Breakthrough with ... \(https://www.quantinuum.com/press-releases/quantinuum-announces-generative-quantum-ai-breakthrough-with-massive-commercial-potential#:~:text=Quantinuum%20Announces%20Generative%20Quantum%20AI,Complex%20Problems%20Impossible%20for\)\)L36](#)]). Sur le plan théorique, une étude dans *Nature* a confirmé la tendance au **gonflement des modèles** : les LLM les plus récents continuent de croître en taille, mais avec des gains de plus en plus marginaux, posant la question de l'approche *scaling* vs *innovations d'architect (AI race in 2025 is tighter than ever before - Nature (https://www.nature.com/articles/d41586-025-01033-y#:~:text=AI%20race%20in%202025%20is,making%20variables%2C%20more%20computing))L27]). Enfin, en **sciences sociales**, des travaux ont évalué l'impact de ChatGPT dans l'éducation : au Royaume-Uni, 92% des étudiants interrogés utilisent des IA comme aide aux dev ([March 2025 round-up of interesting AI news and announcements - Artificial intelligence \(http://nationalcentreforai.jiscinvolve.org/wp/2025/03/27/march-2025-round-up-of-interesting-ai-news-and-announcements#:~:text=%E2%80%98AI%20inclusive%E2%80%99\)\)108](#)], obligeant les universités à repenser l'évaluation (on voit émerger des examens oraux ou des dissertations "assistées" plutôt que bannir l'outil). Une cadre (*AI Literacy Framework*) a même été proposé pour définir les compétences en IA que tout citoyen devrait maîtriser, preuve que l'**acculturation** à l'IA devient un sujet éducatif en ([March 2025 round-up of interesting AI news and announcements - Artificial intelligence \(http://nationalcentreforai.jiscinvolve.org/wp/2025/03/27/march-2025-round-up-of-interesting-ai-news-and-announcements#:~:text=Digital%20Education%20Council%20AI%20Literacy,to%20different%20disciplines%20and%20jurisdictions\)\)L97](#)]). Ces divers résultats montrent que l'IA n'est pas confinée à l'informatique : elle infuse tous les domaines de la recherche, de la chimie à l'histoire, et les chercheurs commencent à la traiter comme un outil standard tout en étudiant ses propres biais et limites (méta-recherche sur les hallucinations, etc.).

Analyse critique : impacts, limites et enjeux

L'effervescence de mars 2025 en IA s'accompagne d'une prise de recul nécessaire sur les **impacts** et les **limites** de ces technologies. **Sur le plan économique et social**, les nouvelles IA promettent des gains de productivité et des services inédits, mais renforcent en parallèle les craintes sur l'emploi et les inégalités. L'autonomisation d'agents comme Manus ou la perspective de robots humanoïdes omniprésents soulèvent la question d'un déplacement, voire d'une disparition, de nombreuses tâches humaines d'ici quelques décenn ([Intelligence Artificielle : Les 5 actualités majeures du 14 mars 2025 \(https://fr.linkedin.com/pulse/intelligence-artificielle-les-5-actualit%C3%A9s-majeures-du-nobili-mpeve#:~:text=2%E2%8F%B8%E2%83%A3%20L%27IA%20et%20la%20robotique,de%20robots%20humano%C3%AFdes%20d%E2%80%99ici%202060\)\)L70](#)]). Les emplois intellectuels ne sont pas épargnés : GPT-4.5 et consorts accomplissent de mieux en mieux des travaux de rédaction, de codage ou d'analyse financière, obligeant les professionnels à **monter en compétences** pour garder une valeur ajoutée (d'où l'urgence de la formation continue, soulignée par de nombreux experts). En même temps, l'histoire montre que chaque révolution technologique crée de nouveaux métiers : on voit émerger des besoins de *spécialistes IA* (prompt engineers, data curators, éthiciens de l'IA), et les pays investissant dans ces compétences pourraient en tirer profit.

Sur le plan de la fiabilité, malgré les progrès, les modèles actuels conservent des **lacunes**. Les phénomènes d'**hallucination** (inventions d'informations factuelles erronées) restent fréquents, même avec GPT-4.5 qui pourtant les réduit en pa ([Introducing GPT-4.5 | OpenAI \(https://openai.com/index/introducing-gpt-4-5#:~:text=We%E2%80%99re%20releasing%20a%20research%20preview,generate%20creative%20insights%20without%20reasoning\)\)148](#)]). Une étude du Tow Center a d'ailleurs mesuré que les chatbots grand public (ChatGPT, Gemini...) fournissaient des réponses incorrectes à plus de 60% de questions de culture générale ([Anthropic adds web search to its Claude chatbot | TechCrunch \(https://techcrunch.com/2025/03/20/anthropic-adds-web-search-to-its-claude-chatbot#:~:text=Claude%E2%80%99s%20ability%20to%20search%20the,with%20the%20reversal%20in%20course\)\)174](#)]). Cela rappelle que ces IA, aussi impressionnantes soient-elles, ne possèdent pas de compréhension infallible du monde. Il en va de même pour les modèles visuels : Midjourney v7 a beau être meilleur, il peut encore produire des

images trompeuses ou des aberrations subtiles. Ces limites techniques posent un enjeu de **confiance** : comment s'assurer qu'une IA utilisée en entreprise, en médecine ou en justice ne commettra pas d'erreur dramatique ? La communauté scientifique travaille sur des solutions (techniques de vérification, d'explicabilité, ou hybrides IA+logique formelle), mais aucune n'est encore totalement satisfaisante à l'heure actuelle.

Un autre **défi** est la **dépendance aux ressources**. La course aux modèles géants et aux services IA mainstream a mis en lumière l'énorme besoin en calcul : OpenAI a dû brider la génération d'images faute de GPU suffisants ([OpenAI says 'our GPUs are melting' as it limits ChatGPT image generation requests | The Verge](https://www.theverge.com/news/637542/chatgpt-says-our-gpus-are-melting-as-it-puts-limit-on-image-generation-requests#:~:text=The%20fervor%20around%20ChatGPT%E2%80%99s%20more,handling%20the%20avalanche%20of%20requests))172) ([https://www.theverge.com/news/637542/chatgpt-says-our-gpus-are-melting-as-it-puts-limit-on-image-generation-requests#:~:text=The%20fervor%20around%20ChatGPT%E2%80%99s%20more,handling%20the%20avalanche%20of%20requests\)\)172](https://www.theverge.com/news/637542/chatgpt-says-our-gpus-are-melting-as-it-puts-limit-on-image-generation-requests#:~:text=The%20fervor%20around%20ChatGPT%E2%80%99s%20more,handling%20the%20avalanche%20of%20requests))172)], et Sam Altman a publiquement quémandé des « lots de 100 000 GPU » pour tenir la char ([OpenAI Shifts Course, Says GPT-5 Coming in 'a Few Months'](https://www.inc.com/ben-sherry/openai-shifts-course-says-gpt-5-coming-in-a-few-months/91172067#:~:text=anything%20the%20company%20has%20delivered,so%20far))117) ([https://www.inc.com/ben-sherry/openai-shifts-course-says-gpt-5-coming-in-a-few-months/91172067#:~:text=anything%20the%20company%20has%20delivered,so%20far\)\)117](https://www.inc.com/ben-sherry/openai-shifts-course-says-gpt-5-coming-in-a-few-months/91172067#:~:text=anything%20the%20company%20has%20delivered,so%20far))117)). Nvidia, de son côté, capitalise sur cette demande explosive en accélérant la cadence des nouveaux GPU haut de gamme ([GTC 2025 – Announcements and Live Updates | NVIDIA Blog](https://blogs.nvidia.com/blog/nvidia-keynote-at-gtc-2025-ai-news-live-updates/?utm_source=chatgpt.com#:~:text=Rubin%20architecture%2C%20designed%20to%20drive))L60) ([https://blogs.nvidia.com/blog/nvidia-keynote-at-gtc-2025-ai-news-live-updates/?utm_source=chatgpt.com#:~:text=Rubin%20architecture%2C%20designed%20to%20drive\)\)L60](https://blogs.nvidia.com/blog/nvidia-keynote-at-gtc-2025-ai-news-live-updates/?utm_source=chatgpt.com#:~:text=Rubin%20architecture%2C%20designed%20to%20drive))L60)). Mais tout le monde n'a pas les moyens d'OpenAI ou de Google : de plus en plus de voix appellent à optimiser l'IA pour la **sobriété**. Cela passe par la recherche de modèles plus petits mais efficaces (d'où l'intérêt des Mistral 3.1 et autres optimisations), et par l'amélioration de l'infrastructure (dissiper moins de chaleur, utiliser du refroidissement plus vert, mutualiser via le cloud pour éviter le gâchis d'équipements sous-utilisés). L'**impact environnemental** de l'IA devient en effet un sujet sensible, d'autant que les grands acteurs restent discrets sur leur consommation énergétique. Des rapports indépendants pointent du doigt les data centers des GAFAM qui s'étendent et engloutissent toujours plus d'électricité (« *Power hungry...* ») comme titrait un article ce mois ([March 2025 round-up of interesting AI news and announcements - Artificial intelligence](http://nationalcentreforai.jiscinvolve.org/wp/2025/03/27/march-2025-round-up-of-interesting-ai-news-and-announcements/#:~:text=Environment))151) ([http://nationalcentreforai.jiscinvolve.org/wp/2025/03/27/march-2025-round-up-of-interesting-ai-news-and-announcements/#:~:text=Environment\)\)151](http://nationalcentreforai.jiscinvolve.org/wp/2025/03/27/march-2025-round-up-of-interesting-ai-news-and-announcements/#:~:text=Environment))151)). À l'heure de la transition écologique, l'IA devra prouver qu'elle peut être une partie de la solution (optimisation énergétique, smart grids) et non une aggravation du problème.

Sur le volet **éthique et juridique**, mars 2025 a montré une nette accélération des efforts de régulation et d'encadrement volontaire. Cependant, un fossé demeure entre la cadence **ultra-rapide de l'innovation** et la lenteur relative des instances décisionnelles. Les législations (telles que l'AI Act européen) ne seront effectives que dans un ou deux ans, or d'ici là les capacités des IA auront encore décuplé. Certains experts craignent une prolifération de faux générés par IA (deepfakes indétectables, textes de propagande à grande échelle) avant que des garde-fous robustes ne soient en place. La coordination internationale reste délicate : les tensions géopolitiques (US vs Chine notamment) compliquent l'adoption de normes communes. En outre, l'**équilibre régulation/innovation** est un exercice d'équilibriste : trop brider pourrait freiner des avancées bénéfiques, mais ne rien faire serait irresponsable. Le code de bonnes pratiques dévoilé cherche cet équilibre ([Intelligence Artificielle : Les 5 actualités majeures du 14 mars 2025](https://fr.linkedin.com/pulse/intelligence-artificielle-les-5-actualit%C3%A9s-majeures-du-nobillmpeve#:~:text=%E2%9A%96%EF%B8%8F%20Face%20%C3%A0%20l%E2%80%99essor%20rapide,transparence%2C%20de%20s%C3%A9curit%C3%A9%20et%20d%E2%80%99impact%20sur%20la%20soci%C3%A9t%C3%A9))133) ([https://fr.linkedin.com/pulse/intelligence-artificielle-les-5-actualit%C3%A9s-majeures-du-nobillmpeve#:~:text=%E2%9A%96%EF%B8%8F%20Face%20%C3%A0%20l%E2%80%99essor%20rapide,transparence%2C%20de%20s%C3%A9curit%C3%A9%20et%20d%E2%80%99impact%20sur%20la%20soci%C3%A9t%C3%A9\)\)133](https://fr.linkedin.com/pulse/intelligence-artificielle-les-5-actualit%C3%A9s-majeures-du-nobillmpeve#:~:text=%E2%9A%96%EF%B8%8F%20Face%20%C3%A0%20l%E2%80%99essor%20rapide,transparence%2C%20de%20s%C3%A9curit%C3%A9%20et%20d%E2%80%99impact%20sur%20la%20soci%C3%A9t%C3%A9))133)), mais son application dépend du bon vouloir des firmes. On observe heureusement que les entreprises majeures commencent à intégrer des **principes éthiques** en interne (OpenAI publie ses valeurs, Anthropic se fonde sur un "Constitutional AI", etc.), souvent sous la pression de l'opinion publique. Le mois de mars a aussi été riche en discussions sur la **transparence** : plusieurs éditeurs (par ex. Wiley) ont émis des directives sur l'utilisation d'IA dans les publications scientifiques ([March 2025 round-up of interesting AI news and announcements - Artificial intelligence](http://nationalcentreforai.jiscinvolve.org/wp/2025/03/27/march-2025-round-up-of-interesting-ai-news-and-announcements/#:~:text=Wiley%20releases%20AI%20guidelines%20for,manuscripts%20while%20using%20AI%20tools))133) ([http://nationalcentreforai.jiscinvolve.org/wp/2025/03/27/march-2025-round-up-of-interesting-ai-news-and-announcements/#:~:text=Wiley%20releases%20AI%20guidelines%20for,manuscripts%20while%20using%20AI%20tools\)\)133](http://nationalcentreforai.jiscinvolve.org/wp/2025/03/27/march-2025-round-up-of-interesting-ai-news-and-announcements/#:~:text=Wiley%20releases%20AI%20guidelines%20for,manuscripts%20while%20using%20AI%20tools))133)), et des médias s'inquiètent de voir leur trafic web cannibalisé par les réponses instantanées de ([March 2025 round-up of interesting AI news and announcements - Artificial intelligence](http://nationalcentreforai.jiscinvolve.org/wp/2025/03/27/march-2025-round-up-of-interesting-ai-news-and-announcements/#:~:text=AI%20search%20summaries%20cannibalise%20academic,of%20research%20more%20clearly%2C%20argues))142) ([http://nationalcentreforai.jiscinvolve.org/wp/2025/03/27/march-2025-round-up-of-interesting-ai-news-and-announcements/#:~:text=AI%20search%20summaries%20cannibalise%20academic,of%20research%20more%20clearly%2C%20argues\)\)142](http://nationalcentreforai.jiscinvolve.org/wp/2025/03/27/march-2025-round-up-of-interesting-ai-news-and-announcements/#:~:text=AI%20search%20summaries%20cannibalise%20academic,of%20research%20more%20clearly%2C%20argues))142)). Cela pose la question de la **créditisation des sources** et de la juste rémunération : si une IA résume un article de presse en donnant la réponse sans que le lecteur ne clique la source, comment le média gagne-t-il sa vie ? Des modèles économiques devront évoluer pour préserver un écosystème de l'information viable à l'ère des réponses directes.

Enfin, un point critique est celui de l'**inclusivité et du biais**. Malgré les efforts d'équilibrage, on constate que les modèles restent biaisés par les données sur lesquelles ils sont entraînés. Les langues peu présentes en corpus restent moins bien servies (d'où l'initiative EuroBert pour combler ce manque en Europe ([Intelligence Artificielle : Les 5 actualités majeures du 14 mars 2025](https://fr.linkedin.com/pulse/intelligence-artificielle-les-5-actualit%C3%A9s-majeures-du-nobillmpeve#:~:text=L%E2%80%99intelligence%20artificielle%20continue%20de%20s%E2%80%99imposer,%C3%A0%20l%E2%80%99information%20dans%20toutes%20les))L81) ([https://fr.linkedin.com/pulse/intelligence-artificielle-les-5-actualit%C3%A9s-majeures-du-nobillmpeve#:~:text=L%E2%80%99intelligence%20artificielle%20continue%20de%20s%E2%80%99imposer,%C3%A0%20l%E2%80%99information%20dans%20toutes%20les\)\)L81](https://fr.linkedin.com/pulse/intelligence-artificielle-les-5-actualit%C3%A9s-majeures-du-nobillmpeve#:~:text=L%E2%80%99intelligence%20artificielle%20continue%20de%20s%E2%80%99imposer,%C3%A0%20l%E2%80%99information%20dans%20toutes%20les))L81))). De même, des biais culturels ou de genre persistent dans les outputs de ChatGPT ou Claude, ce qui peut être problématique si ces outils sont utilisés sans discernement. Mars 2025 a vu quelques polémiques sur des réponses stéréotypées de certaines IA, rappelant la nécessité d'un travail continu sur ce front. La publication du recueil français des 50 termes clés de l'IA (par la ([JFrog s'associe à Hugging Face pour assurer la sécurité des modèles GenAI](https://www.solutions-numeriques.com/jfrog-sassocie-a-hugging-face-pour-assurer-la-s%C3%A9curit%C3%A9-des-mod%C3%A8les-genai/#:~:text=%E2%9E%9C%20IA%20%3A%20le%20recueil,cl%C3%A9s%20fran%C3%A7ais%20compl%C3%A8tement%20pass%C3%A9%20inaper%C3%A7u))141) ([https://www.solutions-numeriques.com/jfrog-sassocie-a-hugging-face-pour-assurer-la-s%C3%A9curit%C3%A9-des-mod%C3%A8les-genai/#:~:text=%E2%9E%9C%20IA%20%3A%20le%20recueil,cl%C3%A9s%20fran%C3%A7ais%20compl%C3%A8tement%20pass%C3%A9%20inaper%C3%A7u\)\)141](https://www.solutions-numeriques.com/jfrog-sassocie-a-hugging-face-pour-assurer-la-s%C3%A9curit%C3%A9-des-mod%C3%A8les-genai/#:~:text=%E2%9E%9C%20IA%20%3A%20le%20recueil,cl%C3%A9s%20fran%C3%A7ais%20compl%C3%A8tement%20pass%C3%A9%20inaper%C3%A7u))141))) vise aussi à diffuser une culture partagée de l'IA pour éviter incompréhensions et fantasmes.

En somme, nous sommes dans une phase où les **opportunités** offertes par l'IA sont immenses – gains de productivité, nouvelles solutions médicales, éducation personnalisée, etc. – mais où les **risques** sont encore mal maîtrisés – erreurs, usage malveillant, concentration du pouvoir, impact sur l'emploi. Mars 2025 a illustré cette dualité avec des avancées technologiques spectaculaires d'un côté, et de l'autre des débuts de réponses (régulation, éthique, outils de sécurisation) encore timides mais indispensables. L'équation à résoudre sera de maximiser les bénéfices tout en minimisant les dommages collatéraux, ce qui requerra une collaboration étroite entre développeurs, utilisateurs, régulateurs et société civile.

Perspectives pour les mois à venir

Les tendances observées en mars laissent entrevoir plusieurs **évolutions clés dans les prochains mois**. D'abord, la course à la **puissance des modèles** va se poursuivre. OpenAI devrait dévoiler **GPT-5** d'ici fin 2025, un modèle annoncé comme unifiant la compréhension non supervisée (type GPT-4.5) et le raisonnement avancé (type o3) dans une seule IA surpuissante ([OpenAI Shifts Course, Says GPT-5 Coming in 'a Few Months'](https://www.inc.com/ben-sherry/openai-shifts-course-says-gpt-5-coming-in-a-few-months/91172067#:~:text=In%20February%2C%20Altman%20posted%20on%205))L83) ([https://www.inc.com/ben-sherry/openai-shifts-course-says-gpt-5-coming-in-a-few-months/91172067#:~:text=In%20February%2C%20Altman%20posted%20on%205\)\)L83](https://www.inc.com/ben-sherry/openai-shifts-course-says-gpt-5-coming-in-a-few-months/91172067#:~:text=In%20February%2C%20Altman%20posted%20on%205))L83)). On peut s'attendre à ce que GPT-5 intègre nativement la **multimodalité totale** (texte, image, son, vidéo) et des fonctionnalités d'*agent* capables d'enchaîner de manière autonome des actions complètes ([OpenAI Shifts Course, Says GPT-5 Coming in 'a Few Months'](https://www.inc.com/ben-sherry/openai-shifts-course-says-gpt-5-coming-in-a-few-months/91172067#:~:text=that%20%E2%80%9Cif%20anyone%20has%20GPU,can%20get%20asap%20please%20call%21%E2%80%9D))) ([https://www.inc.com/ben-sherry/openai-shifts-course-says-gpt-5-coming-in-a-few-months/91172067#:~:text=that%20%E2%80%9Cif%20anyone%20has%20GPU,can%20get%20asap%20please%20call%21%E2%80%9D\)\)](https://www.inc.com/ben-sherry/openai-shifts-course-says-gpt-5-coming-in-a-few-months/91172067#:~:text=that%20%E2%80%9Cif%20anyone%20has%20GPU,can%20get%20asap%20please%20call%21%E2%80%9D)))) ([OpenAI Shifts Course, Says GPT-5 Coming in 'a Few Months'](https://www.inc.com/ben-sherry/openai-shifts-course-says-gpt-5-coming-in-a-few-months/91172067#:~:text=While%20we%20don%E2%80%99t%20know%20exactly,25%2C%20and%20media%20generation))120) ([https://www.inc.com/ben-sherry/openai-shifts-course-says-gpt-5-coming-in-a-few-months/91172067#:~:text=While%20we%20don%E2%80%99t%20know%20exactly,25%2C%20and%20media%20generation\)\)120](https://www.inc.com/ben-sherry/openai-shifts-course-says-gpt-5-coming-in-a-few-months/91172067#:~:text=While%20we%20don%E2%80%99t%20know%20exactly,25%2C%20and%20media%20generation))120))). Sa sortie sera un événement majeur, susceptible de rebattre les cartes si OpenAI tient ses promesses d'améliorations qualitatives. Google de son côté pourrait accélérer l'intégration de Gemini 2.5 dans ses produits :

lors de Google I/O 2025 (attendu en mai), il n'est pas exclu qu'ils annoncent une version **Gemini 3** ou une extension grand public de ces capacités de *reasoning* dans Assistant, Workspace, etc. Anthropic pourrait lancer **Claude 4** (beaucoup de rumeurs autour d'un modèle majeur financé par les 4 milliards d'Ama ([Coud Grok 3 result in Claude 4 and GPT 4.5 to be released earlier...](#))).

(https://www.reddit.com/r/singularity/comments/1ir8nwf/could_grok_3_result_in_claude_4_and_gpt_4_5_to_be_released_earlier...) (<https://blog.promptplayer.com/claude-4/#%7E:text=Claude%204%20Haiku%2C%20Sonnet%2C%20Opus,2025%29L35%20>) , probablement avec un contexte encore plus grand et des "pouvoirs" web accrus. Et n'oublions pas Meta : après LLaMA 3, on chuchote déjà qu'un **LLaMA 4** serait en préparation pour 2025, misant sur une **efficacité paramétrique** plutôt qu'une simple hausse de taille b ([We will get multiple release of Llama 4 in 2025 : r/LocalLLaMA](#) (https://www.reddit.com/r/LocalLLaMA/comments/1hhyozf/we_will_get_multiple_release_of_llama_4_in_2025/#%7E:text=We%20will%20get%20multiple%20release,fit%20into%203%20)).

Parallèlement, la **démocratisation de l'IA** va s'amplifier. La présence de ChatGPT dans les smartphones iOS/Android (via l'API ChatGPT, ou les intégrations Siri suite au partenariat Apple-Op ([OpenAI and Apple announce partnership | OpenAI](#) (<https://openai.com/index/openai-and-apple-announce-partnership/#%7E:text=Apple%20is%20integrating%20ChatGPT%20into,needing%20to%20jump%20between%20tools%29134%20>))) devrait se généraliser, rendant ces assistants omniprésents. Microsoft, qui a intégré GPT-4 dans la suite Office (Microsoft 365 Copilot), va élargir la disponibilité de ces copilotes à tous ses clients professionnels d'ici l'été 2025. On aura alors des **assistants IA dans chaque application de bureau**, changeant la façon de travailler au quotidien. De même, on peut imaginer Google intégrer profondément Gemini dans Gmail, Docs, YouTube (modération et chapitrage automatiques), etc. En somme, l'IA *générative* deviendra moins un produit séparé et de plus en plus une **brique infrastructurelle** invisible mais partout présente, à l'image d'Internet. Cette ubiquité soulève cependant une **attente forte de fiabilité** : la moindre bourde spectaculaire d'un assistant IA intégré (par ex. une erreur dans un email d'entreprise) sera très médiatisée. Les prochains mois seront donc déterminants pour peaufiner ces intégrations et bâtir la confiance des utilisateurs finaux.

Sur le front de la **régulation**, on verra se concrétiser certaines initiatives. L'Union Européenne finalisera sans doute l'AI Act d'ici la fin de sa présidence tournante (mi-2025), ce qui obligera les fournisseurs de modèles à se mettre en conformité (transparence des données, enregistrement auprès d'autorités en fonction du niveau de risque de l'application, etc.). D'autres pays pourraient suivre avec leurs propres lois (les États-Unis travaillent sur un *AI Bill of Rights*, la Chine sur des règles imposant par ex. l'authentification des deepfakes). On assistera probablement aussi à la création de **comités éthiques internes** dans les grandes entreprises utilisatrices d'IA, chargés d'évaluer l'impact des déploiements et de dialoguer avec les régulateurs. La standardisation internationale (ISO) pourrait accoucher de premières normes fin 2025 encadrant par exemple la qualité des jeux de données d'entraînement ou les tests de biais. D'ici là, le **code de bonnes pratiques** lancé en mars servira de guide transitoire, et on verra si les géants respectent réellement leurs engagements volontaires (ce qui influencera la main plus ou moins lourde des régulateurs ensuite).

Techniquement, une **dynamique open source** soutenue va se poursuivre. Mistral AI a annoncé travailler sur un modèle ~30B multimodal d'ici l'automne, qui pourrait constituer une alternative open à GPT-4 en termes d'usages courants. D'autres initiatives communautaires (comme RedPajama, Dolly v3 chez Databricks, etc.) sortiront pour démocratiser l'accès aux LLM. Le rêve d'un « *ChatGPT open source* » de performance équivalente reste vivace – Llama 3.1 s'en approche déjà sur certaines tâches ([Meta releases new Llama 3.1 models, including highly anticipated 405B parameter variant | IBM](#) ([https://www.ibm.com/think/news/meta-releases-llama-3-1-models-405b-parameter-variant/#%7E:text=Gemini%20Ultra%201.0%20\(2853.2\)%20L89%20](https://www.ibm.com/think/news/meta-releases-llama-3-1-models-405b-parameter-variant/#%7E:text=Gemini%20Ultra%201.0%20(2853.2)%20L89%20))) – et il n'est pas impossible qu'en 2025 une solution libre parvienne à combler complètement le fossé. Cela permettrait une plus large adoption encore, mais poserait aussi la question du **contrôle** : un modèle open aussi puissant que GPT-4, accessible à n'importe qui, pourrait tout autant servir des finalités malveillantes (fabrication massive de fake news par des régimes autoritaires, etc.). Ce scénario accroîtrait l'urgence de développer des *contre-mesures* IA (par ex. des détecteurs de texte/image générés) et de mettre en place des **marqueurs d'authenticité** dans les contenus (sujet sur lequel travaillent OpenAI et d'autres).

En **robotique**, les annonces de Nvidia et l'acquisition de Pollen Robotics indiquent qu'on devrait voir plus de synergies IA + robots. Attendez-vous à des démonstrations de robots assistants dans des contextes réels dès la fin d'année 2025 : essais de robots infirmiers dans des cliniques pilotes, robots livreurs autonomes dans certains campus, etc., soutenus par les avancées en vision et en modèles de commandes intelligentes. La question de l'acceptation sociale sera cruciale : ces robots devront gagner la confiance par leur fiabilité et leur conformité aux normes (sécurité, respect de la vie privée s'ils collectent des données). Les prochains mois serviront de test : par petites touches, on verra où ces machines peuvent s'intégrer utilement sans rejet du public. L'autre aspect sera la **standardisation** du "cerveau" des robots : si des modèles comme GR00T N1 (ou d'autres comme *Google Robotics Transformer*) prouvent leur efficacité, ils pourraient devenir des bases communes, accélérant tout le secteur comme ImageNet l'avait fait pour la vision en son temps.

Enfin, il faut s'attendre à de nouvelles **percées scientifiques** grâce à l'IA. La découverte de nouveaux médicaments via des modèles génératifs, par exemple, pourrait connaître un succès retentissant d'ici peu (plusieurs molécules *IA-conçues* entrent en phase d'essais cliniques en 2025). De même en énergie, des IA optimisent déjà la fusion nucléaire expérimentale ou la conception de batteries plus efficaces : une annonce majeure dans l'un de ces domaines attribuée à l'IA créerait un engouement positif. Ces réussites seraient de nature à légitimer davantage l'IA aux yeux du grand public, en montrant des bénéfices tangibles pour la société.

En conclusion, la période à venir s'annonce à la fois **passionnante et critique**. Passionnante car chaque mois apporte son lot d'innovations qui repoussent les limites de ce que les machines peuvent faire ; critique car c'est maintenant que se jouent les **règles du jeu** qui encadreront l'IA pour la prochaine décennie. Mars 2025 nous a donné un aperçu d'un futur très proche où l'intelligence artificielle sera omniprésente, aussi banale qu'Internet, mais il dépend de nous qu'elle soit synonyme de progrès partagé et non de risques incontrôlés.

Les acteurs semblent en avoir conscience : la **coopération** entre industriels, régulateurs, chercheurs et citoyens sera le maître-mot pour **façonner une IA de confiance**. Les prochains mois – et années – verront si nous parvenons collectivement à relever ce défi sans freiner l'élan d'innovation. Une chose est sûre : l'intelligence artificielle continuera d'évoluer à un rythme effréné, et il faudra garder un œil attentif sur chaque nouvelle avancée pour en saisir les implications dans toute leur ampleur.

Sources : Op ([Introducing GPT-4.5 | OpenAI](#) (<https://openai.com/index/introducing-gpt-4-5/#%7E:text=We%E2%80%99re%20releasing%20a%20research%20preview,generate%20creative%20insights%20without%20reasoning%29>)) (OpenAI Shifts Course, Says GPT-5 Coming in 'a Few Months' (<https://www.inc.com/ben-sherry/openai-shifts-course-says-gpt-5-coming-in-a-few-months/91172067/#%7E:text=rolling%20out%20its%20next%20AI,5%2C%20in%20%E2%80%9Ca%20few%20months,%E2%80%9D%29L74%20>)) • Mac4 ([OpenAI dévoile GPT-4.5, plus intelligent et plus "humain"](#) (<https://www.mac4ever.com/ia/187546-openai-devoile-gpt-4-5-plus-intelligent-et-plus-humain/#%7E:text=Avec%20GPT,4%20compris%20les%20attentes%20implicites%29L68%20>)) • Blog du Modéra ([Les 10 modèles d'IA les plus performants en mars 2025](#) (<https://www.blogdumoderateur.com/modeles-ia-plus-performants-mars-2025/#%7E:text=conserve%20sa%20premi%C3%A8re%20position%20acquise,cinq%C3%A8me%20position%2C%20tandis%20qu%E2%80%99o1%20est%29>)) (Mistral AI lance Small 3.1, un modèle léger qui prétend surpasser la concurrence (<https://www.blogdumoderateur.com/mistral-ai-lance-small-3-1-surpasser-concurrence/#%7E:text=Cette%20nouvelle%20version%2C%20qui%20%E2%80%99appuie,offrant%20une%20vitesse%20d%E2%80%99inf%C3%A9rence%20inf%C3%A9rieure%29>)) • OpenAI (Inc. N ([OpenAI Shifts Course, Says GPT-5 Coming in 'a Few Months'](#) (<https://www.inc.com/ben-sherry/openai-shifts-course-says-gpt-5-coming-in-a-few-months/91172067/#%7E:text=In%20February%2C%20Altman%20posted%20on,5%29>)) (OpenAI Shifts Course, Says GPT-5 Coming in 'a Few Months' (<https://www.inc.com/ben-sherry/openai-shifts-course-says-gpt-5-coming-in-a-few-months/91172067/#%7E:text=The%20capacity%20concerns%20come%20on,%E2%80%9D%29106%20>)) • Googl ([Gemini 2.5: Our](#)

newest Gemini model with thinking. (https://blog.google/technology/google-deepmind/gemini-model-thinking-updates-march-2025/#:~:text=Today%20we%E2%80%99re%20introducing%20Gemini%202_LMArena%20by%20a%20significant%20margin)) (Gemini 2.5: Our newest Gemini model with thinking. (<https://blog.google/technology/google-deepmind/gemini-model-thinking-updates-march-2025/#:~:text=Introducing%20Gemini%202>))357】 • TechCr (Anthropic adds web search to its Claude chatbot | TechCrunch (<https://techcrunch.com/2025/03/20/anthropic-adds-web-search-to-its-claude-chatbot/#:~:text=that%20had%20long%20eluded%20it>)) (Anthropic adds web search to its Claude chatbot | TechCrunch (https://techcrunch.com/2025/03/20/anthropic-adds-web-search-to-its-claude-chatbot/#:~:text=Claude%E2%80%99s%20ability%20to%20search%20the_with%20the%20reversal%20in%20course))173】 • Anthr (Claude takes research to new places | Anthropic (<https://www.anthropic.com/news/research#:~:text=Google%20Workspace>)) (Anthropic adds web search to its Claude chatbot | TechCrunch (https://techcrunch.com/2025/03/20/anthropic-adds-web-search-to-its-claude-chatbot/#:~:text=Web%20search%20is%20available%20now_sites%20to%20inform%20certain%20responses))148】 • Nvidia (GTC 2025 – Announcements and Live Updates | NVIDIA Blog. (https://blogs.nvidia.com/blog/nvidia-keynote-at-gtc-2025-ai-news-live-updates/?utm_source=chatgpt.com#:~:text=In%20a%20video%2C%20Huang%20announced_generalized%20humanoid%20reasoning%20and%20skills)) (GTC 2025 – Announcements and Live Updates | NVIDIA Blog. (https://blogs.nvidia.com/blog/nvidia-keynote-at-gtc-2025-ai-news-live-updates/?utm_source=chatgpt.com#:~:text=Rubin%20architecture%2C%20designed%20to%20drive))L60】 • Upskil (IA : 3 grandes avancées en mars 2025 qui vont tout changer - Upskilling. (https://upskilling.com/ia-3-grandes-avancees-en-mars-2025-gui-vont-tout-change/#:~:text=Le%206%20mars%202025%2C%20la_besoin%20d%E2%80%99une%20intervention%20humaine%20constante)) (IA : 3 grandes avancées en mars 2025 qui vont tout changer - Upskilling. (https://upskilling.com/ia-3-grandes-avancees-en-mars-2025-gui-vont-tout-change/#:~:text=La%20mont%C3%A9e%20en%20puissance%20des_IA%20plus%20transparente%20et%20responsable))L75】 • LinkedIn (Intelligence Artificielle : Les 5 actualités majeures du 14 mars 2025. (https://fr.linkedin.com/pulse/intelligence-artificielle-les-5-actualit%C3%A9s-majeures-du-nobili-mpeve#:~:text=E2%9A%96%EF%B8%8F%20Face%20%C3%A0%20l%E2%80%99essor%20rapide_transparence%2C%20de%20s%C3%A9curit%C3%A9%20et%20d%E2%80%99)) (Intelligence Artificielle : Les 5 actualités majeures du 14 mars 2025. (https://fr.linkedin.com/pulse/intelligence-artificielle-les-5-actualit%C3%A9s-majeures-du-nobili-mpeve#:~:text=Les%20grandes%20institutions%20financi%C3%A8res%20acc%C3%A9l%C3%A8rent_bulles%20sp%C3%A9culatives%20aliment%C3%A9es%20par%20l%E2%80%99)) • Solutions (JFrog s'associe à Hugging Face pour assurer la sécurité des modèles GenAI (https://www.solutions-numeriques.com/jfrog-sassocie-a-hugging-face-pour-assurer-la-securite-des-modeles-genai/#:~:text=Cette%20int%C3%A9gration%20consistera%2C%20sch%C3%A9matiquement%2C%20%C3%A0_Karas%2C%20CTO%20de%20JFrog%20Security))133】 • TechCr (Stability AI's new AI model turns photos into 3D scenes | TechCrunch (https://techcrunch.com/2025/03/18/stability-ais-new-ai-model-turns-photos-into-3d-scenes/#:~:text=Stability%20AI%20has%20released%20a_with%20realistic%20depth%20and%20perspective)) (Runway releases an impressive new video-generating AI model | TechCrunch (https://techcrunch.com/2025/03/31/runway-releases-an-impressive-new-video-generating-ai-model/#:~:text=Called%20Gen_perspectives%20and%20positions%20within%20scenes))149】 • Blogdumodéra (Test de Midjourney v7 : un modèle bourré de qualités... mais aussi de défauts (<https://www.blogdumoderateur.com/test-midjourney-v7/#:~:text=styles%20graphiques>))L51】 • IBM (Meta releases new Llama 3.1 models, including highly anticipated 405B parameter variant | IBM (<https://www.ibm.com/think/news/meta-releases-llama-3-1-models-405b-parameter-variant#:~:text=Gemini%20Ultra%201.0%202853.2>))L89】 • Jisc AI Ce (March 2025 round-up of interesting AI news and announcements - Artificial intelligence (http://nationalcentreforai.jiscinvolve.org/wp/2025/03/27/march-2025-round-up-of-interesting-ai-news-and-announcements/#:~:text=Introducing%20NextGenAI%3A%20A%20consortium%20to_using%20AI%20to%20accelerate%20research))-L4】 , etc. (voir liens ci-dessus).