

Taux d'hallucination dans les systemes d'IA contemporains (2023-2025)

Sources verifiees uniquement : etudes evaluees par des pairs, references officielles et rapports de recherche institutionnelle | Prepare aux fins d'evaluation institutionnelle | Toutes les sources verifiees en avril 2026

50-83 %

Taux d'hallucination dans une etude sur des vignettes cliniques portant sur 6 grands modeles (Mount Sinai / Nature Medicine, 2025)

17-34 %

Taux d'hallucination dans des outils juridiques RAG specialises commercialises comme 'sans hallucination' (Stanford HAI, 2024)

33-51 %

Taux d'hallucination des modeles de raisonnement o3 d'OpenAI sur des references factuelles ouvertes (SimpleQA / PersonQA, 2025)

Constat central de l'ensemble des donnees examinees : L'hallucination n'est pas une limite temporaire que l'amelioration des modeles permettra d'eliminer. Il s'agit d'une propriete structurelle persistante des systemes d'IA a inference prioritaire, dont la forme evolue a mesure que les modeles progressent. Sur les taches contraintes et ancrees, les taux diminuent. Sur le raisonnement complexe, les questions en domaine ouvert et les entrees adversariales, les taux demeurent eleves - et dans certains cas augmentent a mesure que les modeles deviennent plus performants.

BASE DE DONNEES PROBANTES - ETUDE PAR ETUDE

Omar et al. (2025) - Attaques adversariales par hallucination en aide a la decision clinique

CLINIQUE / NATURE MEDICINE

MODELES TESTES

GPT-4o, Distilled-DeepSeek-Llama, Phi-4, Gemma-2-27B-it, Qwen-2.5-72B, Llama-3.3-70B (6 modeles, 5 400 sorties)

DEFINITION ET METHODE

300 vignettes validees par des medecins, chacune contenant un detail clinique fabrique (valeur de laboratoire, signe ou affection). Hallucination = le modele elabore sur le detail fabrique comme s'il etait reel.

TAUX RAPPORTES

Parametre par default global : 65,9 %
DeepSeek-Llama : 80-83 %
GPT-4o (meilleur cas) : 50-53 %
Avec invite d'attenuation : 44,2 %
GPT-4o + attenuation : 23 %

Evalue par des pairs dans Communications Medicine (Nature Publishing Group), PMID : PMC12318031 | pmc.ncbi.nlm.nih.gov/articles/PMC12318031 | Source verifiee en avril 2026

Magesh et al. (2024) - Hallucination dans les outils juridiques d'IA RAG specialises

JURIDIQUE / STANFORD HAI

MODELES TESTES

Lexis+ AI (LexisNexis), Westlaw AI-Assisted Research (Thomson Reuters), Ask Practical Law AI (Thomson Reuters), GPT-4

DEFINITION ET METHODE

Plus de 200 requetes juridiques ouvertes construites manuellement. Hallucination = enonce juridique incorrect ou citation mal ancree ne soutenant pas l'affirmation formulee.

TAUX RAPPORTES

Westlaw AI-Assisted Research : >34 %
Lexis+ AI : >17 %
Ask Practical Law AI : >17 %
GPT-4 usage general (etude anterieure) : 58-82 %

Stanford RegLab et HAI | hai.stanford.edu/news/ai-trial-legal-models-hallucinate-1-out-6-or-more-benchmarking-queries | Source verifiee en avril 2026
Constat de : meme les outils juridiques RAG specialises commercialises comme 'sans hallucination' ont produit des informations incorrectes dans plus d'une requete sur six. La RAG reduit les erreurs mais ne les elimine pas.

Graffius (2026) - Les hallucinations de l'IA s'améliorent-elles ou s'aggravent-elles?

MULTIMODELE / SYNTHÈSE DE
RÉFÉRENCES

MODELES TESTES

OpenAI o3, o1; Gemini-2.0-Flash; plusieurs fournisseurs. Références : Vectara HHEM, SimpleQA, PersonQA, évaluations domaine-spécifiques.

DEFINITION ET METHODE

Synthèse des données de références 2024-2025. Deux tendances divergentes : tâches de synthèse ancrées (taux en baisse) vs raisonnement complexe en domaine ouvert (taux en hausse).

TAUX RAPPORTES

Synthèse ancrée, meilleurs modèles : 0,7-1,5 %
OpenAI o3 sur PersonQA / SimpleQA : 33-51 %
OpenAI o1 (génération précédente) : ~16 %
Évaluations de tâches larges 2025 : 3-20 %+

DOI : 10.13140/RG.2.2.33179.53285 | scottgraffius.com/blog/files/ai-hallucinations-2026.html | Source vérifiée en avril 2026

Constat clé : les modèles de raisonnement de nouvelle génération hallucinent davantage sur les tâches factuelles complexes que les modèles antérieurs. Une capacité de raisonnement accrue ne se traduit pas automatiquement par une meilleure fiabilité factuelle.

Chelli et al. (2024) - Taux d'hallucination dans les citations académiques générées par l'IA

CLINIQUE / JMIR

MODELES TESTES

GPT-3.5 (ChatGPT), GPT-4, Bard/Gemini

DEFINITION ET METHODE

Référence générée ne correspondant à aucune métadonnée de publication réelle. Revue systématique des citations générées par l'IA pour la littérature médicale.

TAUX RAPPORTES

GPT-3.5 : 39,6 %
GPT-4 : 28,6 %
Bard/Gemini : 91,4 %

Journal of Medical Internet Research, 2024, e53164 | jmir.org/2024/1/e53164

Tableau de classement Vectara HHEM (2025) - Référence de synthèse ancrée

REFERENCE

MODELES TESTES

Gemini-2.0-Flash, GPT-4o, GPT-o1, DeepSeek-R1, DeepSeek-V3 et autres

DEFINITION ET METHODE

Hughes Hallucination Evaluation Model (HHEM) : fidélité de la sortie du modèle par rapport à un document source fourni. Tâche de synthèse ancrée et contrainte - conditions optimales.

TAUX RAPPORTES

Gemini-2.0-Flash : ~0,7 %
GPT-o1 / GPT-4o : ~0,9-1,9 %
DeepSeek-R1 : 14,3 %
DeepSeek-V3 : 3,5 %

Tableau de classement Vectara | github.com/vectara/hallucination-leaderboard

Remarque : il s'agit de taux optimaux sur une tâche de synthèse contrainte - non représentatifs des performances en domaine ouvert ou en raisonnement complexe.

CE QUE LES DONNÉES ÉTABLISSENT

Le dossier de recherche couvrant 2023-2025 établit trois constats directement pertinents pour la gouvernance institutionnelle de l'IA.

Premièrement, les taux d'hallucination dépendent de la tâche, non du modèle.

La même famille de modèles qui atteint moins de 1 % sur une tâche contrainte de synthèse documentaire peut produire des erreurs dans 33 à 51 % des requêtes factuelles en domaine ouvert. Le tableau de classement Vectara montre que les meilleurs modèles approchent 0,7 % sur la synthèse ancrée. La fiche système d'OpenAI pour o3 documente des taux d'erreur de 33 à 51 % sur SimpleQA et PersonQA. Ces résultats ne sont pas contradictoires - ils reflètent la même architecture sous-jacente fonctionnant dans des conditions différentes. Les institutions qui évaluent l'IA sur la base de références fournies par les fournisseurs observent généralement des performances optimales dans des conditions contraintes, et non les performances sur les types de questions complexes et spécifiques aux politiques que leurs utilisateurs poseront réellement.

Deuxièmement, la RAG réduit l'hallucination mais ne l'élimine pas.

L'évaluation indépendante par Stanford HAI des principales plateformes d'IA juridique - des outils commercialisés comme 'sans hallucination' par LexisNexis et Thomson Reuters - a révélé des taux d'hallucination de 17 à 34 % dans des conditions réelles de

requetes juridiques. Il s'agit de systemes RAG specialises avec des bases de donnees juridiques organisees. Le constat n'est pas que la RAG echoue - elle reduit substantiellement les erreurs par rapport aux modeles polyvalents - mais qu'elle n'offre aucune garantie architecturale. Les taux d'hallucination dependent de la qualite de la recuperation, de la complexite des requetes et de la presence de fausses premisses, elements qui varient dans les environnements institutionnels reels.

Troisiemement, les modeles plus performants peuvent davantage halluciner sur certaines taches.

La synthese Graffius (2026) documente un phenomene contre-intuitif : la serie o3 d'OpenAI, optimisee pour le raisonnement complexe, hallucine a 33-51 % sur la memorisation factuelle en domaine ouvert - plus du double du taux de la serie o1 anterieure, etabli a environ 16 %. L'etude clinique du Mount Sinai a constate que l'attenuation par invite reduisait le taux d'hallucination adversarial de GPT-4o de 53 % a 23 % - une amelioration significative, mais representant encore un taux d'echec de 23 % sous la meilleure strategie d'attenuation disponible. L'hallucination n'est pas un probleme de qualite de modele que la prochaine generation resoudra. C'est une propriete structurelle des architectures a inference prioritaire.

IMPLICATION POUR LA GOUVERNANCE

Pour les institutions reglementees, la question pertinente n'est pas de savoir si les systemes d'IA actuels hallucinent a des taux acceptables sur des references controlees. Plusieurs le font. La question est de savoir si les conditions dans lesquelles des taux acceptables sont atteints peuvent etre maintenues de maniere fiable dans un deploiement institutionnel - et quelles sont les consequences lorsqu'elles ne le sont pas.

L'etude juridique de Stanford HAI illustre le mieux cet ecart. Les plateformes qui se reclamaient 'sans hallucination' sur la base de metriques etroites d'exactitude des citations ont produit des informations juridiques incorrectes dans plus d'une requete sur six dans des conditions reelles. Une universite, un hopital ou un organisme gouvernemental qui s'appuie sur les donnees de reference des fournisseurs pour satisfaire ses obligations de gouvernance se fonde sur des chiffres de performance optimaux qui peuvent ne pas tenir dans les conditions reelles d'utilisation institutionnelle.

Une architecture qui conditionne l'existence de l'inference a une autorisation prealable a l'execution ne rivalise pas avec ces systemes sur le taux d'hallucination. Elle supprime les conditions structurelles dans lesquelles l'hallucination se produit pour les requetes non autorisees ou hors portee - non pas en atteignant un taux plus bas, mais en empechant l'inference de s'executer lorsqu'aucune source faisant autorite n'existe pour la soutenir.

REFERENCES VERIFIEES

1. Omar M. et al. (2025). Analyse multi-modele montrant la vulnerabilite elevee des grands modeles de langage aux attaques adversariales par hallucination en aide a la decision clinique. Communications Medicine (Nature). PMID : PMC12318031. pmc.ncbi.nlm.nih.gov/articles/PMC12318031
2. Magesh V. et al. (2024). Sans hallucination? Evaluation de la fiabilite des principaux outils de recherche juridique par IA. Stanford RegLab / HAI. hai.stanford.edu/news/ai-trial-legal-models-hallucinate-1-out-6-or-more-benchmarking-queries
3. Graffius S.M. (2026). Les hallucinations de l'IA s'ameliorent-elles ou s'aggravent-elles? Nous avons analyse les donnees. DOI : 10.13140/RG.2.2.33179.53285. scottgraffius.com/blog/files/ai-hallucinations-2026.html
4. Chelli M. et al. (2024). Taux d'hallucination et exactitude des references de ChatGPT et Bard pour les revues systematiques. Journal of Medical Internet Research, e53164. jmir.org/2024/1/e53164
5. Vectara (2025). Tableau de classement du Hughes Hallucination Evaluation Model (HHEM). github.com/vectara/hallucination-leaderboard
6. OpenAI (2025). Fiche systeme o3 et o4-mini. Documente des taux d'hallucination de 33-51 % sur les references PersonQA et SimpleQA. Reference et synthetise dans Graffius (2026).