

Membership Inference Attacks LLM 2026 : Vol



16 mai
2026



Mis à jour le 17 mai
2026



20 min de
lecture



3491
mots



MIA permettent de prouver qu'un texte était dans le training set d'un LLM. 0.65-0.85 sur Llama 4, implications RGPD majeures.

À RETENIR

A retenir — Membership Inference Attacks LLM

MIA détermine si un texte était dans le training set d'un LLM. AUC observée sur Llama 4 70B, 0.74 sur Claude 4.5 (Min-K%).

Familles d'attaque : *loss-based* (Shi et al., 2023), *Min-K%* (Shi et al., 2023) *Model*, *Zlib ratio*.

Implications RGPD : un sujet peut prouver que ses données personnelles ont été dans le modèle — droit à l'effacement applicable.

Defenses : **DP-SGD** (Abadi et al., 2016), réduction de la précision (deduplication), gradient clipping, fine-tuning sans LoRA-Mer

In projet cybersécurité
Réponse sous 24h

Devis
gratuit



Litiges 2026 : *plaintiffs* utilisent MIA pour fonder leurs actions contre OpenAI, Stability AI, Meta. Charge probatoire elevee mais MIA y contribue.

Les **membership inference attacks (MIA)** sur LLM repondent a une question simple mais juridiquement explosive : "Ce texte etait-il dans le training set du modele ?". En 2022, Shih et al. publient *Detecting Pretraining Data from Large Language Models* et demontrent que, par une analyse statistique sur les log-probabilities, on distingue les samples membres (in-training) des non-membres avec un AUC > 0.7 sur GPT-3, Llama, OPT. Trois ans plus tard, les MIA sont au centre de litiges juridiques (NYT vs OpenAI, Getty vs Stability) et de la conformite RGPD (droit a l'effacement). Cet article presente la mecanique mathematique, le code Python, les benchmarks et des defenses (DP-SGD, deduplication). Pour les controleurs RGPD et les juristes specialises, cela representent l'argumentaire technique le plus solide pour fonder le droit a l'effacement. Cet article approfondit la **membership inference** sur Llama 4, Claude 4.5, GPT-5 et autres, avec implications concretes pour les fine-tunings d'entreprise.

1. Genese et etat de l'art

Le concept de MIA est introduit par Shokri et al. (2017) *Membership Inference Attacks on Machine Learning Models*, sur des classifieurs d'images. Pour les LLM, la litterature recente (2022-2023) :

Carlini et al. (2022) *The Privacy Risks of Memorization* — demonstration de memorisation parfaite sur GPT-2.

Carlini et al. (2023) *Quantifying Memorization Across Neural Language Models* — etude de la memorisation.

Réponse au 24h Shi et al. (2023) *Detecting Pretraining Data from Large Language Models*

Devis
gratuit



WIKIMIA,

Réponse sous 24h

Devis
gratuit →