



# LLM Model Extraction 2026 : Voler un GPT

📅 16 mai 2026 • 🔄 Mis à jour le 17 mai 2026 • ⌚ 20 min de lecture • ☰ 3757 mots



Carlini et al. (2024) extraient la matrice de projection finale de GPT-3.5 pour l'extraction par distillation atteint 92% de fidelity pour 50000 USD.

## À RETENIR

### A retenir — LLM Model Extraction

**Carlini et al. (2024)** extraient la matrice de projection embeddings → logits **200 USD** en queries API.

Black-box distillation Llama 4 vers proxy 7B : **91% accuracy** retention sur M 50000 USD, 200M tokens.

Defenses 2026 : *logit truncation* (top-k visible), *differential privacy* sur output Kirchenbauer.

In projet cybersécurité  
Réponse sous 24h  
Impact business : vol de propriété intellectuelle, modèles concurrents, contournement de licences spécialisées

Devis gratuit →

AI Act article 53 : les fournisseurs GPAI doivent documenter les contre-mesures  
Litiges en cours US (NYT vs OpenAI, Getty vs Stability).

Le **model extraction llm**, ou *model stealing*, est l'attaque qui menace directement les propriétaires des laboratoires d'IA. L'idée : un attaquant qui ne connaît pas les poids d'un LLM peut les extraire à partir d'interrogations API. Trois familles existent : (1) **fonctionnal extraction** via les capacités spécifiques (fine-tuning specialise). En 2026, le cout d'extraction d'un LLM est estimé à (50k-500k USD selon la fidelity attendue), mais le ROI attaquant peut être énorme pour un gouvernement adverse, ou un acteur cherchant à republier le modèle sans licence. Pour les fournisseurs de LLM propriétaires, la documentation des contre-mesures (AI Act a une exigence réglementaire incontournable, doublée d'une réalité compétitive où ByteDance et plusieurs startups asiatiques ont prouvé la viabilité de l'extraction industrielle. Cet article explore le **model extraction llm** sur les frontières 2026.

## 1. Genese et etat de l'art

L'attaque de model extraction remonte à Tramer et al. (2016) *Stealing Machine Learning Models from Prediction APIs*, applique à l'époque aux SVM et logistic regression. Pour les NN, les travaux de Carlini et al. (2019) *Model Extraction via API* démontrent la possibilité de cloner via queries crafted. Pour les LLM, les premiers travaux sont :

**Wallace et al. (2020)** : *Imitation Attacks and Defenses for Black-box Machine Learning Models* - extraction de modèles de traduction par la rétention sur Google Translate avec ~100k queries.

**Krishna et al. (2020)** : extraction de BERT-base finetuned sur des données random, etc.

Réponse sous 24h

Devis  
gratuit →

---

Réponse sous 24h

Devis  
gratuit →