



# Multi-Turn Jailbreaks 2026 : Crescendo & Ske



16 mai  
2026



Mis à jour le 17 mai  
2026



20 min de  
lecture



3690  
mots



Les jailbreaks multi-tour (Crescendo, Skeleton Key, Many-Shot) exploitent le contexte des LLM modernes. ASR 80-95%, defenses limitées.

## À RETENIR

### A retenir — Multi-Turn Jailbreaks

**Crescendo** (Russinovich et al., Microsoft 2024) construit une escalade en 5 tours. ASR **87%** sur GPT-5.

**Skeleton Key** (Russinovich, 2024) annule l'alignement par injection de meta-prompt. ASR **92%** sur GPT-4o avant patch.

**Many-Shot Jailbreaking** (Anil et al., Anthropic 2024) exploite le long-contexte (256 tokens) : 256 exemples in-context defont l'alignement.

L'alignement single-turn ne couvre **pas** les attaques multi-tours. Réponse sous 24h. **Devis gratuit** → sur N tours. classifieur de session (cumulative risk score).

Sur agents autonomes (LangChain, CrewAI), un jailbreak multi-turn déclenche des actions outils — impact business critique.

Le **multi-turn jailbreak llm** est aujourd'hui le vecteur d'attaque le plus sous-estimé. Pendant que la littérature académique se concentre sur les attaques single-turn (Cross-Contextual, adversarial), trois familles d'attaques conversationnelles (*Crescendo*, *Skeleton Key*, *Jailbreaking*) atteignent 80 à 95% de taux de succès sur les LLM 2024. La raison : l'alignement RLHF est entraîné sur des dialogues majoritairement courts, et les fenêtres de contexte modernes (200k-2M tokens) ouvrent un espace d'attaque que les classifieurs single-turn ne voient pas. Cet article présente la mécanique, le code Python d'exploitation, et les défenses réellement en production. Pour les RSSI déployant un chatbot client ou un agent autonome, ces attaques sont la priorité réglementaire et opérationnelle. L'[AI Act](#) article 15 impose la documentation des tests multi-turn.

## 1. Genèse et état de l'art

L'idée que la conversation longue affaiblit l'alignement remonte aux observations empiriques sur ChatGPT (jailbreaks "DAN" en plusieurs étapes, role-play prolongé). La première étude académique vient d'Anthropic en avril 2024 : Anil et al. publient *Many-shot Jailbreaking* montrant qu'en injectant 256 exemples de Q&A "non alignés" en in-context, on dégrade systématiquement l'alignement de Claude, GPT-4, Gemini.

En août 2024, Mark Russinovich (Microsoft) publie deux articles fondateurs :

**Crescendo**, une attaque graduelle qui démarre par une question benigne et évolue vers la requête malveillante.

Réponse sous 24h

Devis  
gratuit



---

Réponse sous 24h

Devis  
gratuit →