

ANTHROPIC

The Permanent Defender Advantage has a Single Point of Failure

AI, Attackers, and the Confidential Computing



Jason Clinton
Deputy CISO

Computation used to train notable artificial intelligence systems, by domain

Computation is measured in total petaFLOP, which is 10^{15} floating-point operations¹. Estimated from AI literature, albeit with some uncertainty. Estimates are expected to be accurate within a factor of 2, or a factor of 5 for recent undisclosed models like GPT-4.

Training computation (petaFLOP) (plotted on a logarithmic axis)



Data source: Epoch AI (2025)

OurWorldinData.org/artificial-intelligence | CC BY

Note: The Executive Order on AI refers to a directive issued by President Biden on October 30, 2023, aimed at establishing guidelines and standards for the responsible development and use of artificial intelligence within the United States.

1. Floating-point operation A floating-point operation (FLOP) is a type of computer operation. One FLOP represents a single arithmetic operation involving floating-point numbers, such as addition, subtraction, multiplication, or division.

A view from where we sit

**The agentic era is here.
For attackers too.**

BY MATTHEW GAULT SECURITY JUN 4, 2025 6:00 AM

The Rise of 'Vibe Hacking' Is the Next AI Nightmare

In the very near future, victory will belong to the savvy blackhat hacker who uses AI to generate code at scale.

SEPTEMBER 15, 2025

Vibe hacking: How AI-driven cybercrime outpaces EDR and signature defenses

EXCLUSIVE ARTIFICIAL INTELLIGENCE [Follow](#)

Chinese Hackers Used Anthropic's AI to Automate Cyberattacks

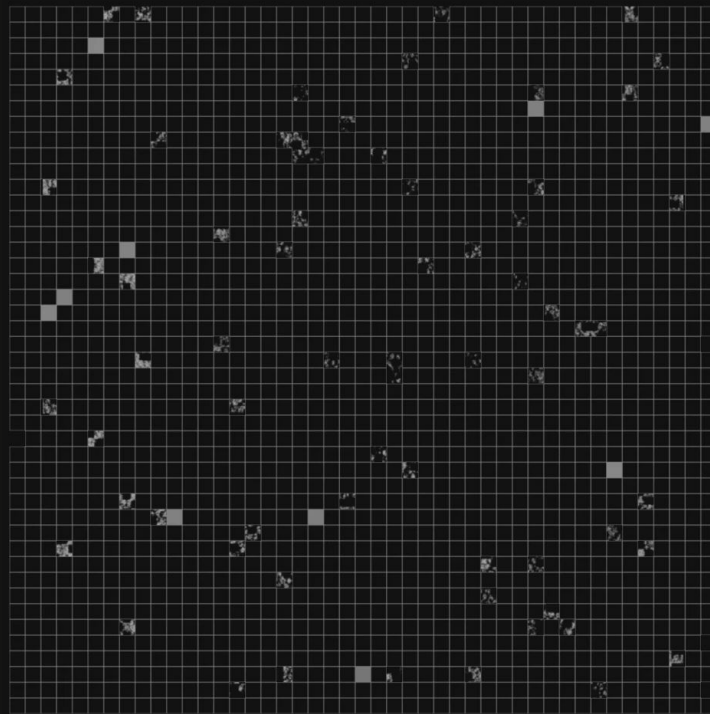
The use of AI automation in hacks is a growing trend that gives hackers additional scale and speed

ANTHROPIC

Project Glasswing

AI now finds what humans missed for decades

AI models have reached a level of coding capability where they can surpass all but the most skilled humans at finding and exploiting software vulnerabilities.



The permanent defender advantage is coming



Software with vulnerabilities
can no longer be released

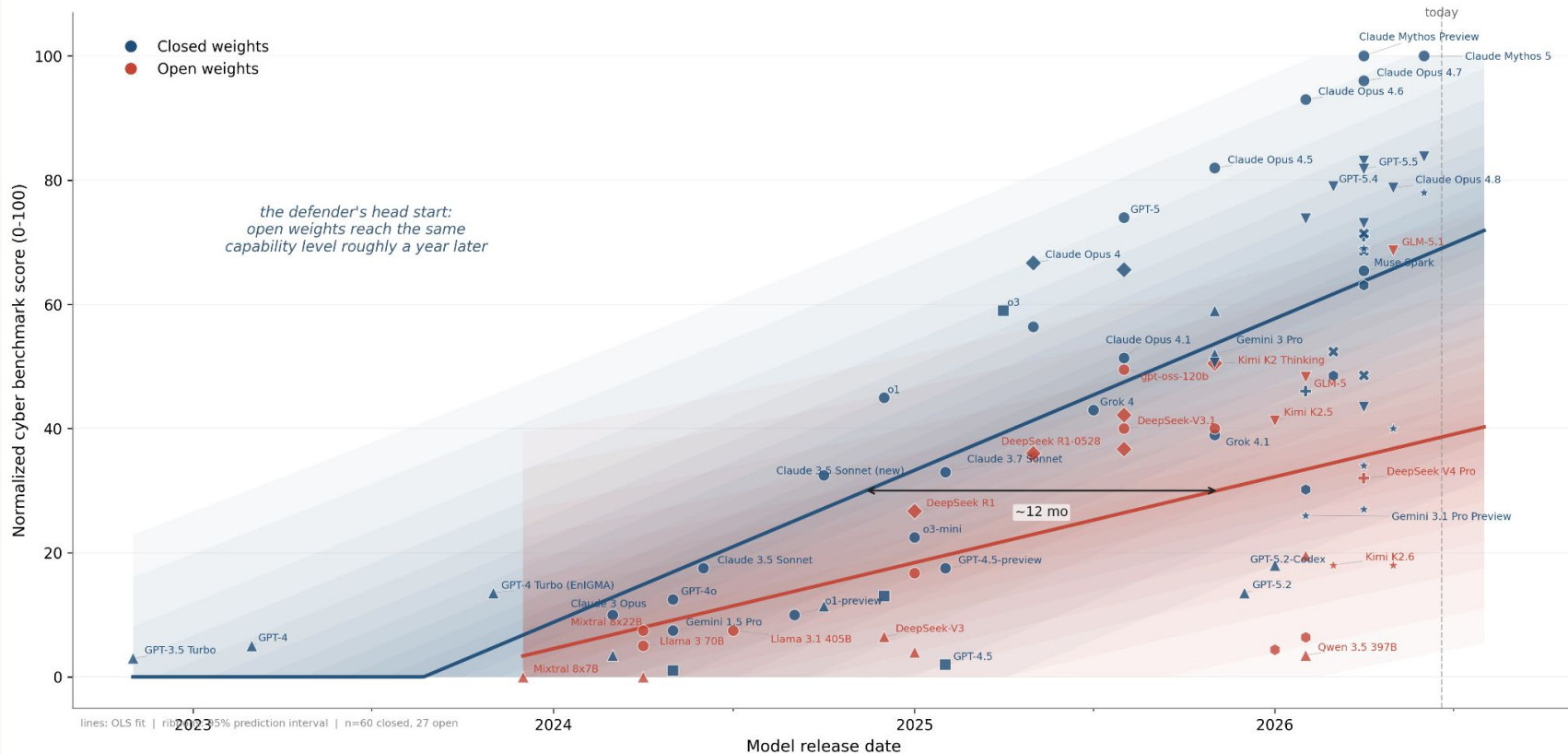


There's only so many ways for
software to be wrong in a way
that is useful for an attacker



It will take longer for OT and
critical infrastructure

Offensive cyber capability over time: closed vs open weights (SOTA-at-release models only)



- Benchmark (marker shape)
- Cybench (40 pro CTFs)
 - ◆ CVE-Bench (NIST CAISI)
 - ▽ CyberGym (1,507 OSS-Fuzz vulns; Berkeley)
 - OpenAI CTF suite, professional tier
 - ⊕ CTF-Archive-Diamond (NIST CAISI, 285 hard)
 - ⊙ SEC-bench Pro (V8/Firefox/kernel PoC)
 - ▲ NYU CTF Bench (200 CSAW challenges)
 - ⊗ UK AISI Advanced Cyber, Expert tier
 - ★ ExploitBench (CMU V8 16-rung ladder, Cap%)

Preserving the defender advantage

How we defend for the next 7-10 months

Defender playbook



Design for breach (zero-trust)

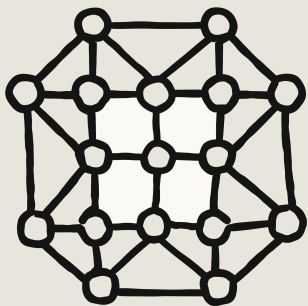


D&R with scaling laws AI



Vulnerability remediation





Don't leak the model weights
while making AI as widely
available to defenders as possible.

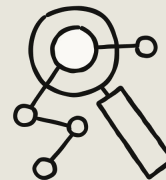
Confidential computing's contribution



Advances in every hardware platform in use today:
Trainium, TPU, and NVIDIA



Enables containment while accelerating reach



Every industry player now agrees that CoCo is part of the solution

ANTHROPIC