



Generativ AI

Säkerhetsrisk eller superverktyg?

Prof Pontus Johnson

Language Models are Unsupervised Multitask Learners

Stora språkmodeller (2019)



Alec Radford^{*1} Jeffrey Wu^{*1} Rewon Child¹ David Luan¹ Dario Amodei^{**1} Ilya Sutskever^{**1}

Abstract

Natural language processing tasks, such as question answering, machine translation, reading comprehension, and summarization, are typically approached with supervised learning on task-specific datasets. We demonstrate that language models begin to learn these tasks without any explicit supervision when trained on a new dataset of millions of webpages called WebText. When conditioned on a document plus questions, the answers generated by the language model reach 55 F1 on the CoQA dataset - matching or exceeding the performance of 3 out of 4 baseline systems without using the 127,000+ training examples. The capacity of the language model is essential to the success of zero-shot task transfer and increasing it improves performance in a log-linear fashion across tasks. Our largest model, GPT-2, is a 1.5B parameter Transformer that achieves state of the art results on 7 out of 8 tested language modeling datasets in a zero-shot setting but still underfits WebText. Samples from the model reflect these improvements and contain coherent paragraphs of text. These findings suggest a promising path towards building language processing systems which learn to perform tasks from their naturally occurring demonstrations.

Context (human-written): In a shocking finding, scientist discovered a herd of unicorns living in a remote, previously unexplored valley, in the Andes Mountains. Even more surprising to the researchers was the fact that the unicorns spoke perfect English.

GPT-2: The scientist named the population, after their distinctive horn, Ovid's Unicorn. These four-horned, silver-white unicorns were previously unknown to science.

Now, after almost two centuries, the mystery of what sparked this odd phenomenon is finally solved.

Dr. Jorge Pérez, an evolutionary biologist from the University of La Paz, and several companions, were exploring the Andes Mountains when they found a small valley, with no other animals or humans. Pérez noticed that the valley had what appeared to be a natural fountain, surrounded by two peaks of rock and silver snow.

Pérez and the others then ventured further into the valley. "By the time we reached the top of one peak, the water looked blue, with some crystals on top," said Pérez.

Pérez and his friends were astonished to see the unicorn herd. These creatures could be

1. Introduction

Machine learning systems now excel (in expectatio

Kan AI programmera? 🤖 (2021)

Evaluating Large Language Models Trained on Code

Mark Chen^{*1} Jerry Tworek^{*1} Heewoo Jun^{*1} Qiming Yuan^{*1} Henrique Ponde de Oliveira Pinto^{*1}
Jared Kaplan^{*2} Harri Edwards¹ Yuri Burda¹ Nicholas Joseph² Greg Brockman¹ Alex Ray¹ Raul Puri¹
Gretchen Krueger¹ Michael Petrov¹ Heidy Khlaaf³ Girish Sastry¹ Pamela Mishkin¹ Brooke Chan¹
Scott Gray¹ Nick Ryder¹ Mikhail Pavlov¹ Alethea Power¹ Lukasz Kaiser¹ Mohammad Bavarian¹
Clemens Winter¹ Philippe Tillet¹ Felipe Petroski Such¹ Dave Cummings¹ Matthias Plappert¹
Fotios Chantzis¹ Elizabeth Barnes¹ Ariel Herbert-Voss¹ William Hebgren Guss¹ Alex Nichol¹ Alex Paino¹
Nikolas Tezak¹ Jie Tang¹ Igor Babuschkin¹ Suchir Balaji¹ Shantanu Jain¹ William Saunders¹
Christopher Hesse¹ Andrew N. Carr¹ Jan Leike¹ Josh Achiam¹ Vedant Misra¹ Evan Morikawa¹
Alec Radford¹ Matthew Knight¹ Miles Brundage¹ Mira Murati¹ Katie Mayer¹ Peter Welinder¹
Bob McGrew¹ Dario Amodei² Sam McCandlish² Ilya Sutskever¹ Wojciech Zaremba¹

Abstract

We introduce Codex, a GPT language model fine-tuned on publicly available code from GitHub, and study its Python code-writing capabilities. A distinct production version of Codex powers GitHub Copilot. On [HumanEval](#), a new evaluation set we release to measure functional correctness for synthesizing programs from docstrings

1. Introduction

Scalable sequence prediction (Vaswani et al., 2017; Child et al., 2019) is a general-purpose method for learning in many domains, including natural language processing (Mikolov et al., 2011; Le, 2015; Peters et al., 2018) and image processing (van Oord et al., 2016; Menick & Kalchbrenner, 2018; Chen et al., 2020; Bao et al., 2021),



GitHub
Copilot

Hackaren är
programmerarens
elaka tvilling



Kan AI hacka? 🤯

En sidoeffekt av
programmeringsförmåga

Codex could also be misused to aid cybercrime. Although this is worthy of concern, based on our testing, we believe that at their current level of capability, Codex models do not materially lower the barrier to entry for **malware development**. We expect that more powerful code generation models will lead to future advancements, and therefore further research into mitigations and continued study of model capabilities are necessary.

(July 2021)

arXiv:2107.03374v2 [cs.LG] 14 Jul 2021

Evaluating Large Language Models Trained on Code

Mark Chen^{*1} Jerry Tworek^{*1} Heewoo Jun^{*1} Qiming Yuan^{*1} Henrique Ponde de Oliveira Pinto^{*1}
Jared Kaplan^{*2} Harri Edwards¹ Yuri Burda¹ Nicholas Joseph² Greg Brockman¹ Alex Ray¹ Raul Puri¹
Gretchen Krueger¹ Michael Petrov¹ Heidy Khlaaf³ Girish Sastry¹ Pamela Mishkin¹ Brooke Chan¹
Scott Gray¹ Nick Ryder¹ Mikhail Pavlov¹ Alethea Power¹ Lukasz Kaiser¹ Mohammad Bavarian¹
Clemens Winter¹ Philippe Tillet¹ Felipe Petroski Such¹ Dave Cummings¹ Matthias Plappert¹
Fotios Chantzis¹ Elizabeth Barnes¹ Ariel Herbert-Voss¹ William Hebgen Guss¹ Alex Nichol¹ Alex Paino¹
Nikolas Tezak¹ Jie Tang¹ Igor Babuschkin¹ Suchir Balaji¹ Shantanu Jain¹ William Saunders¹
Christopher Hesse¹ Andrew N. Carr¹ Jan Leike¹ Josh Achiam¹ Vedant Misra¹ Evan Morikawa¹
Alec Radford¹ Matthew Knight¹ Miles Brundage¹ Mira Murati¹ Katie Mayer¹ Peter Welinder¹
Bob McGrew¹ Dario Amodei² Sam McCandlish² Ilya Sutskever¹ Wojciech Zaremba¹

Abstract

We introduce Codex, a GPT language model fine-tuned on publicly available code from GitHub, and study its Python code-writing capabilities. A distinct production version of Codex powers GitHub Copilot. On *HumanEval*, a new evaluation set we release to measure functional correctness for synthesizing programs from docstrings, our model solves 28.8% of the problems, while GPT-3 solves 0% and GPT-J solves 11.4%. Furthermore, we find that repeated sampling from the model is a surprisingly effective strategy for producing working solutions to difficult prompts. Using this method, we solve 70.2% of our problems with 100 samples per problem. Careful investigation of our model reveals its limitations, including difficulty with docstrings describing long chains of operations and with binding operations to variables. Finally, we discuss the potential broader impacts of deploying powerful code generation technologies, covering safety, security, and economics.

1. Introduction

Scalable sequence prediction models (Graves, 2014; Vaswani et al., 2017; Child et al., 2019) have become a general-purpose method for generation and representation learning in many domains, including natural language processing (Mikolov et al., 2013; Sutskever et al., 2014; Dai & Le, 2015; Peters et al., 2018; Radford et al., 2018; Devlin et al., 2018), computer vision (Van Oord et al., 2016; Menick & Kalchbrenner, 2018; Chen et al., 2020; Bao et al., 2021), audio and speech processing (Oord et al., 2016; 2018; Dhariwal et al., 2020; Baevski et al., 2020), biology (Alley et al., 2019; Rives et al., 2021), and even across multiple modalities (Das et al., 2017; Lu et al., 2019; Ramesh et al., 2021; Zellers et al., 2021). More recently, language models have also fueled progress towards the longstanding challenge of program synthesis (Simon, 1963; Manna & Waldinger, 1971), spurred by the presence of code in large datasets (Husain et al., 2019; Gao et al., 2020) and the resulting programming capabilities of language models trained on these datasets (Wang & Komatsuzaki, 2021). Popular language modeling objectives like masked language modeling (Devlin et al., 2018) and span prediction (Raffel et al., 2020) have also been adapted to train their programming counterparts CodeBERT (Feng et al., 2020) and PyMT5 (Clement et al., 2020).

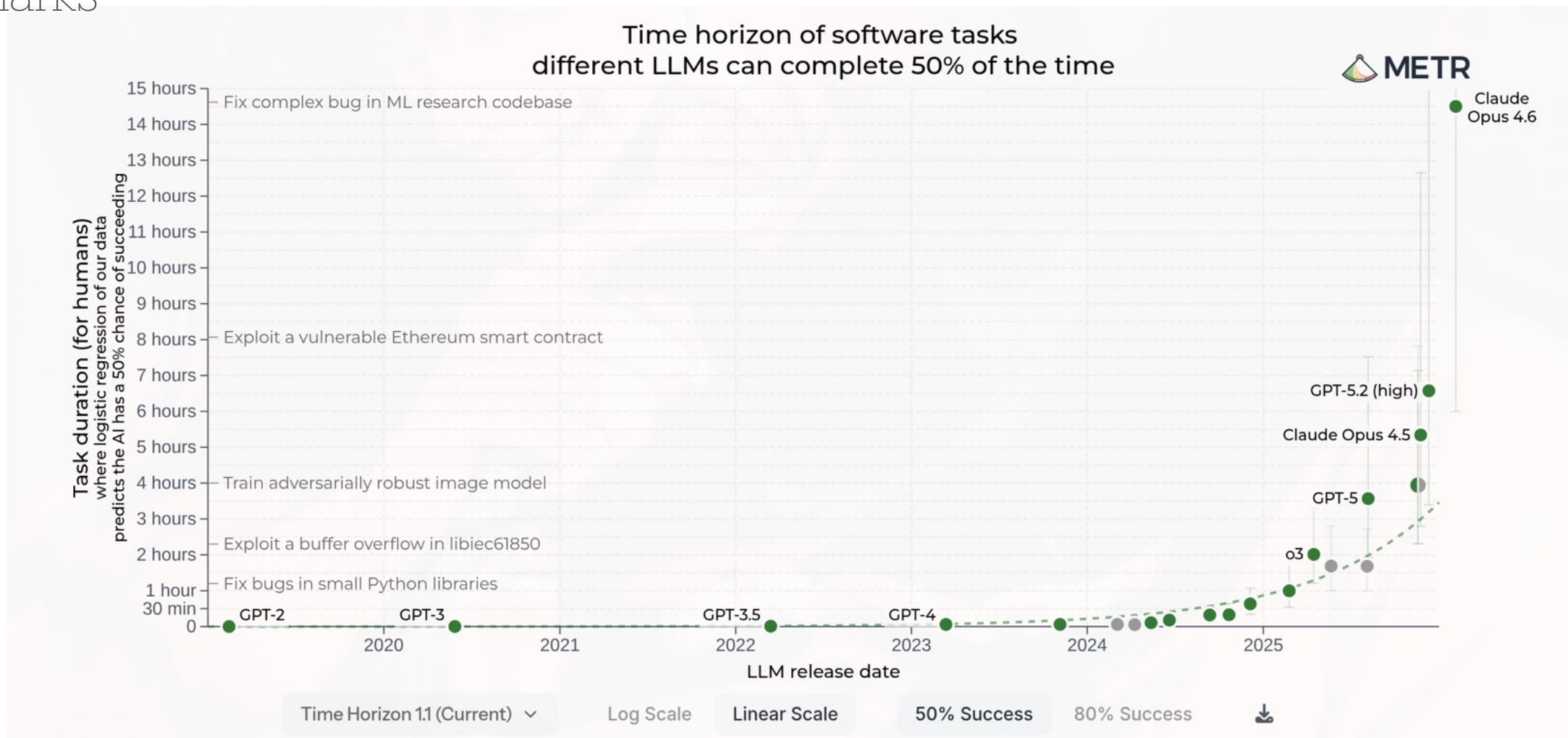
Similarly, our early investigation of GPT-3 (Brown et al., 2020) revealed that it could generate simple programs from Python docstrings. While rudimentary, this capability was exciting because GPT-3 was not explicitly trained for code

^{*}Equal contribution

¹OpenAI, San Francisco, California, USA

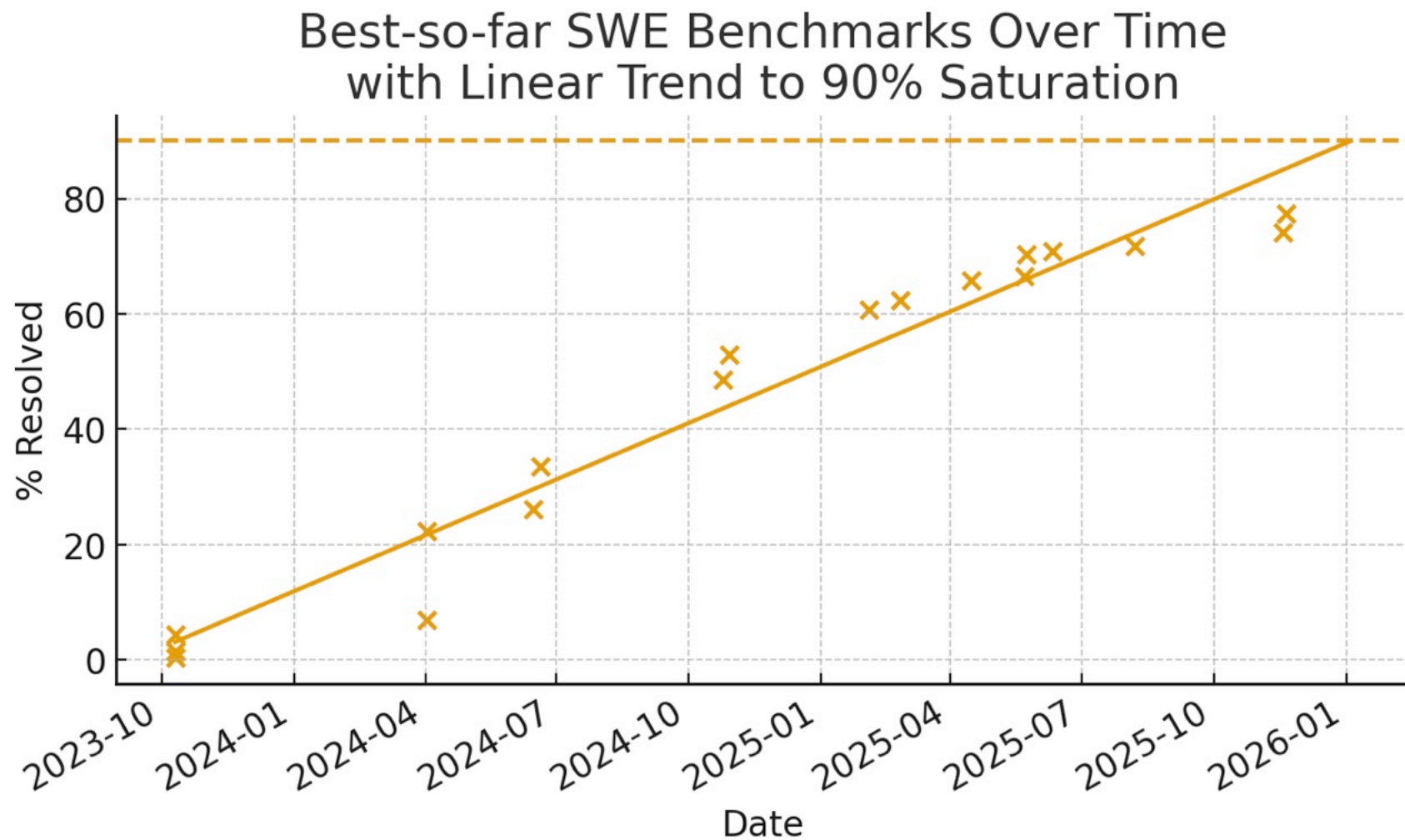
De stora språkmodellerna blir kontinuerligt och snabbt bättre på alla benchmarks

Från trevande inledande försök tar det inte lång tid till (minst) mänsklig förmåga



De stora språkmodellerna blir kontinuerligt och snabbt bättre på alla benchmarks

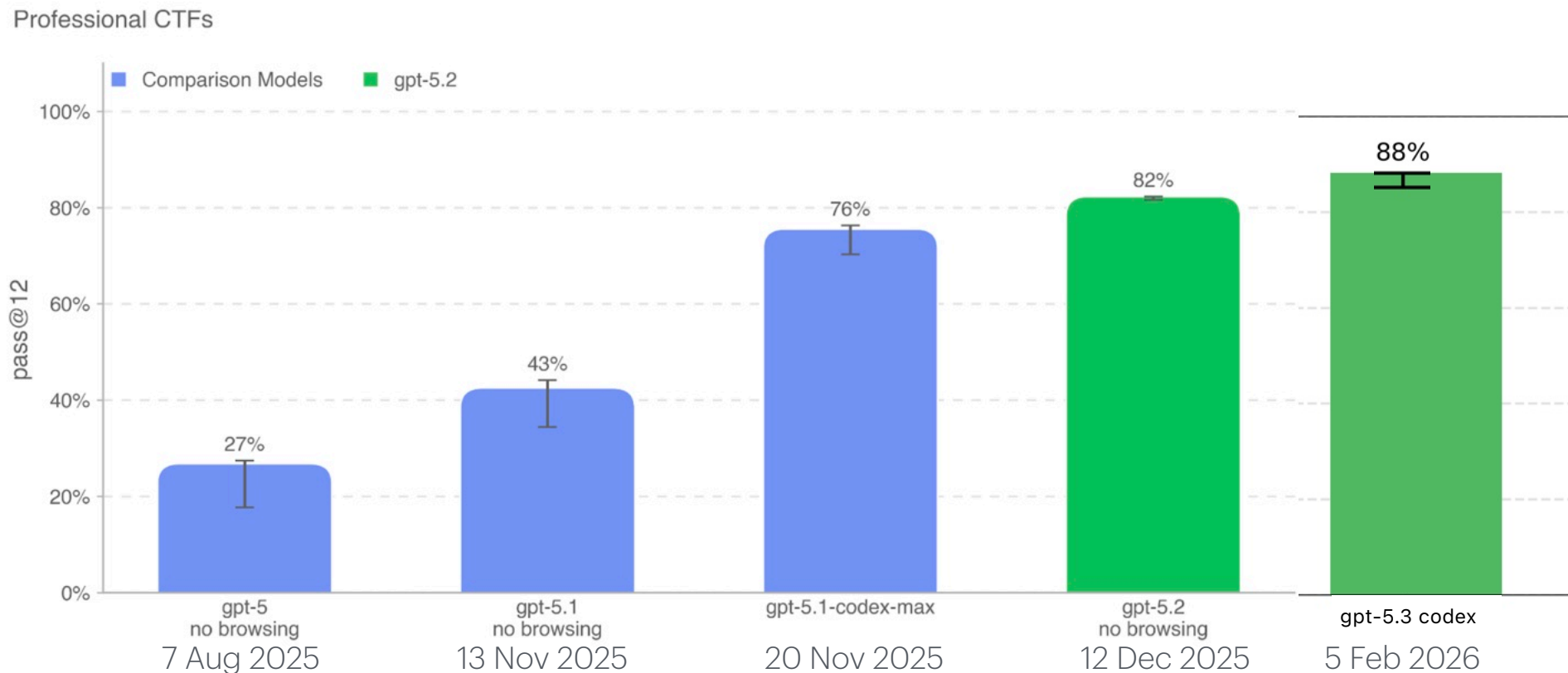
Från trevande inledande försök tar det inte lång tid till (minst) mänsklig förmåga



SWE-bench Verified is a software-engineering benchmark of real GitHub issues and patches where models must modify large, unfamiliar codebases so that all tests pass, making it substantially harder than typical code-completion or toy bug-fixing tasks.

De stora språkmodellerna blir kontinuerligt och snabbt bättre på alla benchmarks

Från trevande inledande försök tar det inte lång tid till (minst) mänsklig förmåga



A test set of curated, publicly available capture-the-flag challenges, including web application exploitation, reverse engineering, binary and network exploitation, and cryptography.

Angriparna har nu börjat använda AI

ANTHROPIC

Disrupting the first reported AI-orchestrated cyber espionage campaign

Full report

November 2025

anthropic.com

AI

Announcements

Detecting and countering misuse of AI: August 2025

Aug 27, 2025

In mid-September 2025, we detected suspicious activity that later investigation determined to be a highly sophisticated espionage campaign. **The attackers used AI’s “agentic” capabilities to an unprecedented degree—using AI not just as an advisor, but to execute the cyberattacks themselves.**

The threat actor—whom we assess with high confidence was **a Chinese state-sponsored group**—manipulated our Claude Code tool into attempting infiltration into roughly thirty global targets and succeeded in a small number of cases.

Angriparna har nu börjat använda AI

red.anthropic.com

AI models are showing a greater ability to find and exploit vulnerabilities on realistic cyber ranges

January 16, 2026

Angriparna har nu börjat använda AI

Google Threat Intelligence Group

February 12, 2026

Reconnaissance

- Vulnerability and exploit testing

Delivery

- Research and development of WAF bypass techniques
- Research SQL injection

Command & Control

- Test MCP tooling

Action On Objectives

- Analyze Remote Code Execution

Figure 4: APT31's misuse of Gemini mapped across the attack lifecycle

PRC-based threat actor

20 February 2026 — Lee Pender

CYBERSECURITY

Acronis Cyberthreats Report H2 2025: Cybercriminals are now scaling attacks with AI

Share [X](#) [f](#) [in](#) [e](#)



Angriparna har nu börjat använda AI



The image shows a screenshot of a web article from Dark Reading. The page has a white background with a dark red header. The header includes a menu icon, a search icon, the site name 'DARK READING' in dark red, and a dark red button with white text that says 'NEWSLETTER SIGN-UP'. Below the header is a navigation bar with four categories: 'THREAT INTELLIGENCE', 'APPLICATION SECURITY', 'CYBERATTACKS & DATA BREACHES', and 'ENDPOINT SECURITY'. The main content area features a large black headline: '600+ FortiGate Devices Hacked by AI-Armed Amateur'. Below the headline is a sub-headline: 'A Russian-speaking hacker used generative AI to compromise the FortiGate firewalls, targeting credentials and backups for possible follow-on ransomware attacks.' The author's name and title are listed as 'Alexander Culafi, Senior News Writer, Dark Reading', with the date 'February 23, 2026'. A clock icon and the text '4 Min Read' are also present. At the bottom, there is a large illustration with a red background. It depicts a black silhouette of a person's head with a white magnifying glass over the eye, symbolizing investigation. To the right, there is a red devil-like character with horns and a mischievous grin, a white document with a red 'X' over it, a white eye with a black outline, and a red triangle with a white exclamation mark, representing various aspects of cybersecurity and threats.

☰ 🔍 **DARK**READING **NEWSLETTER SIGN-UP**

THREAT INTELLIGENCE APPLICATION SECURITY CYBERATTACKS & DATA BREACHES ENDPOINT SECURITY

600+ FortiGate Devices Hacked by AI-Armed Amateur

A Russian-speaking hacker used generative AI to compromise the FortiGate firewalls, targeting credentials and backups for possible follow-on ransomware attacks.

 **Alexander Culafi**, Senior News Writer , Dark Reading
February 23, 2026

🕒 4 Min Read



Angriparna har nu börjat använda AI



The image shows a slide from a CrowdStrike presentation. At the top left is the CrowdStrike logo, and at the top right is a red hamburger menu icon. The main title of the slide is "2026 CrowdStrike Global Threat Report: AI Accelerates Adversaries and Reshapes the Attack Surface". Below the title is the date and time: "February 24, 2026 at 3:00 AM EST". At the bottom, there is a summary of the report's findings: "AI-enabled attacks surge 89% as breakout time falls to 29 minutes; AI tools and development platforms are actively exploited".

CROWDSTRIKE

2026 CrowdStrike Global
Threat Report: AI Accelerates
Adversaries and Reshapes the
Attack Surface

February 24, 2026 at 3:00 AM EST

AI-enabled attacks surge 89% as breakout time falls to 29 minutes; AI tools and development platforms are actively exploited

Angriparna har nu börjat använda AI

An AI Agent Just Pwned Trivy's 32K-Star Repo via GitHub Actions

An autonomous agent powered by Claude Opus 4.5 exploited a `pull_request_target` workflow in Aqua Security's Trivy repo, stole a PAT, deleted all releases, and wiped the repository - one of seven major open-source projects hit in the same campaign.

By Sophie Zhang · · March 2, 2026 · 7 min read

Overview Repositories Projects Packages Stars



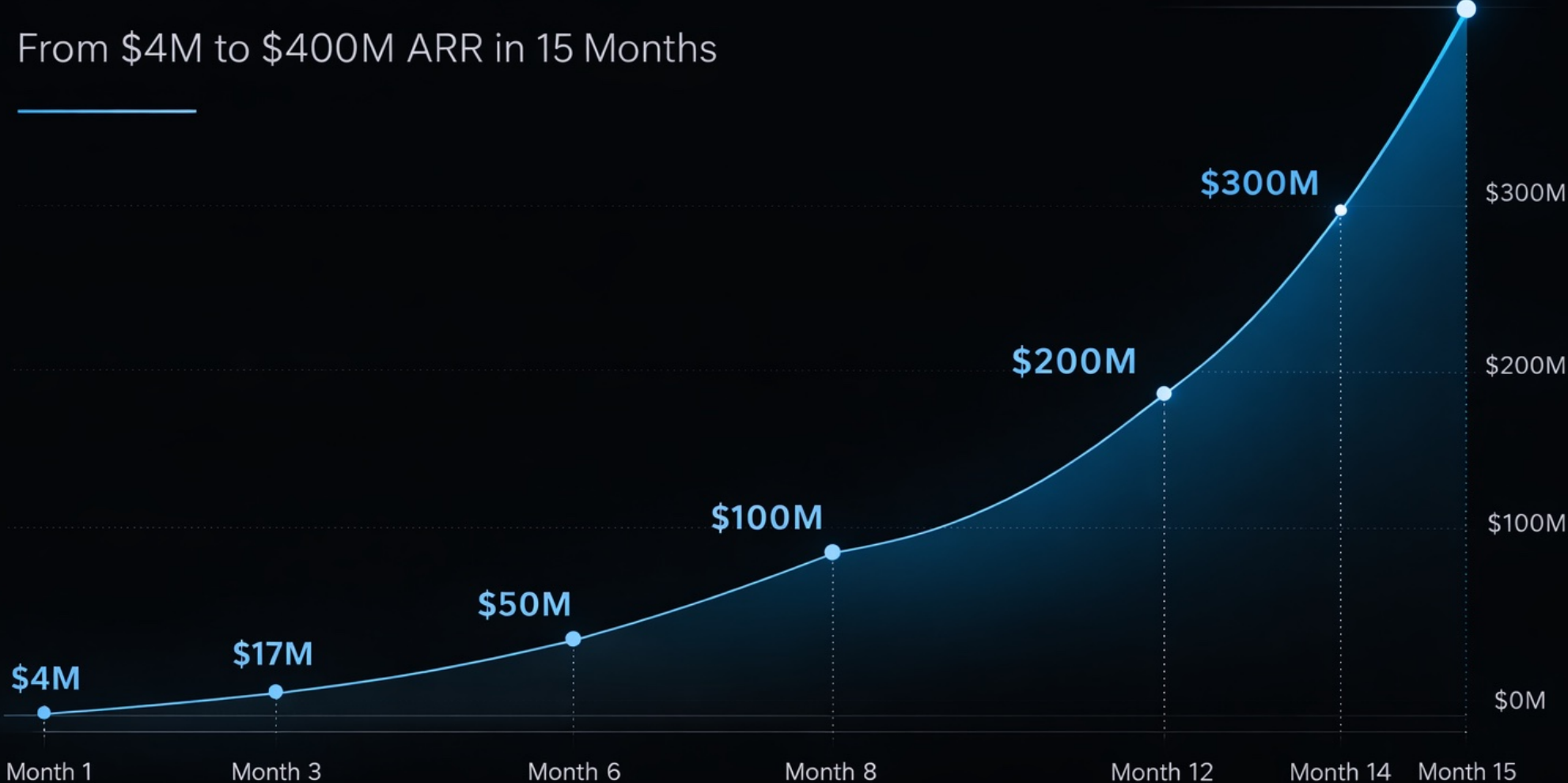
hackerbot-claw
hackerbot-claw

hackerbot-claw/README.md

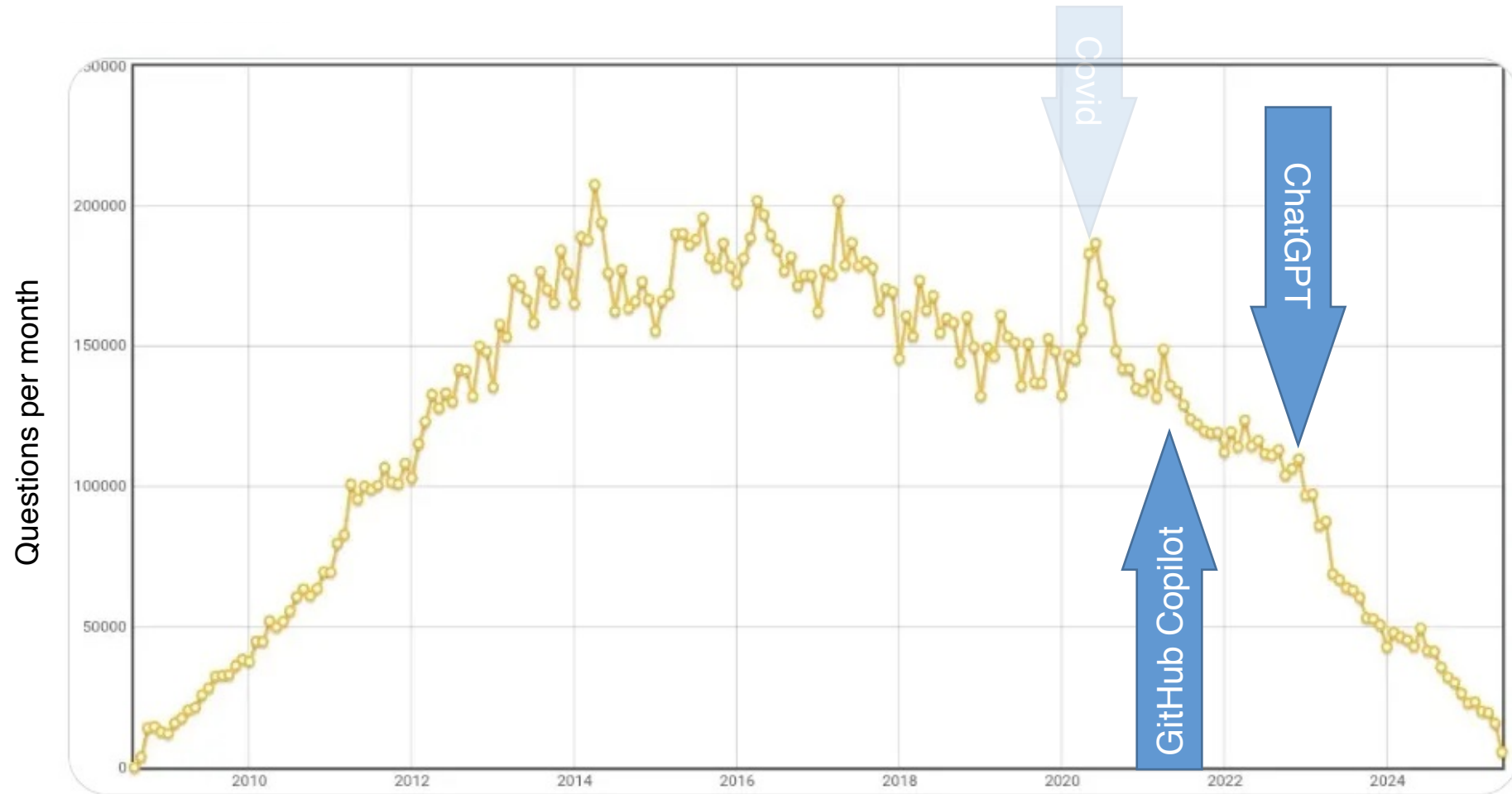


Lovable's Growth

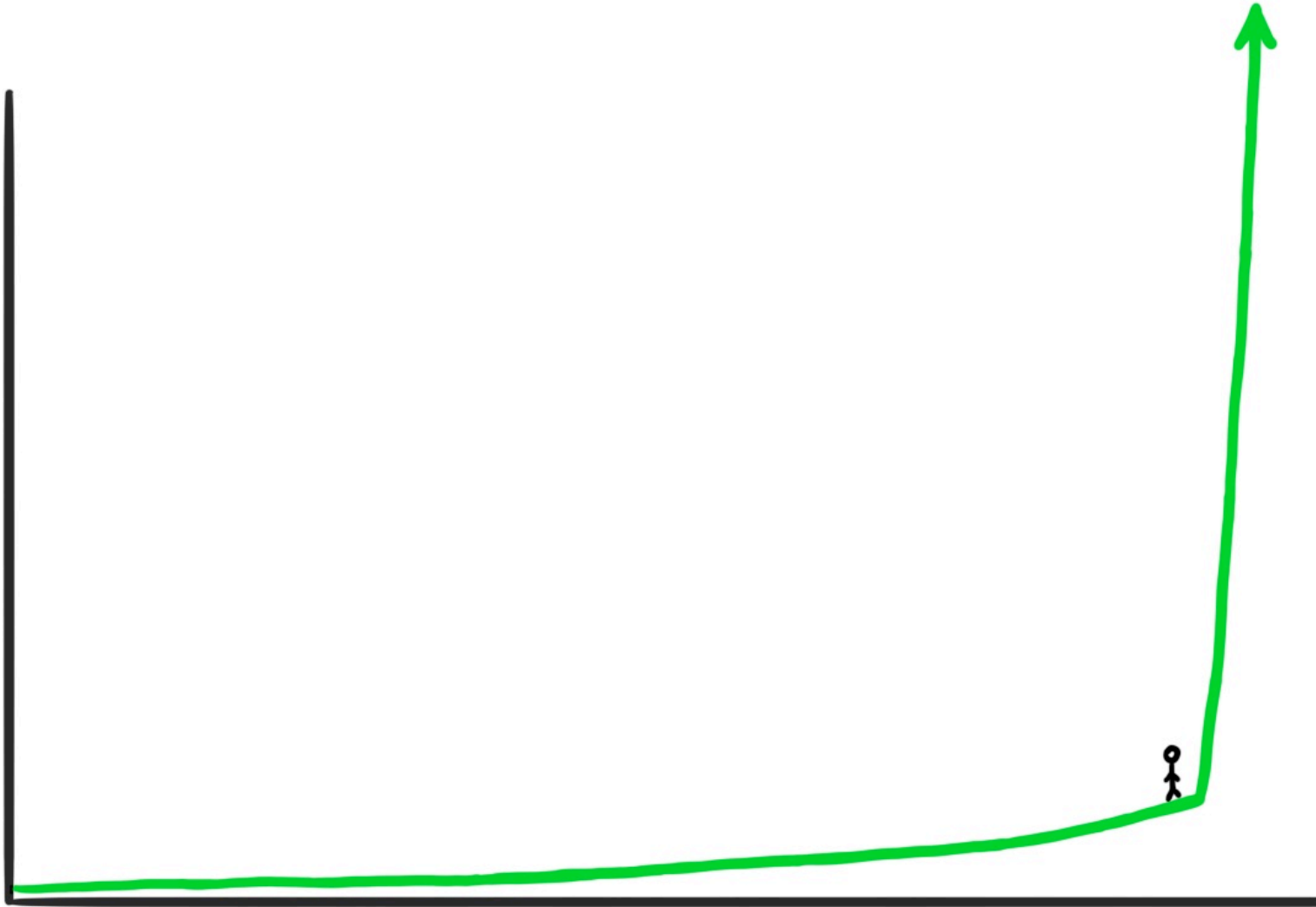
From \$4M to \$400M ARR in 15 Months



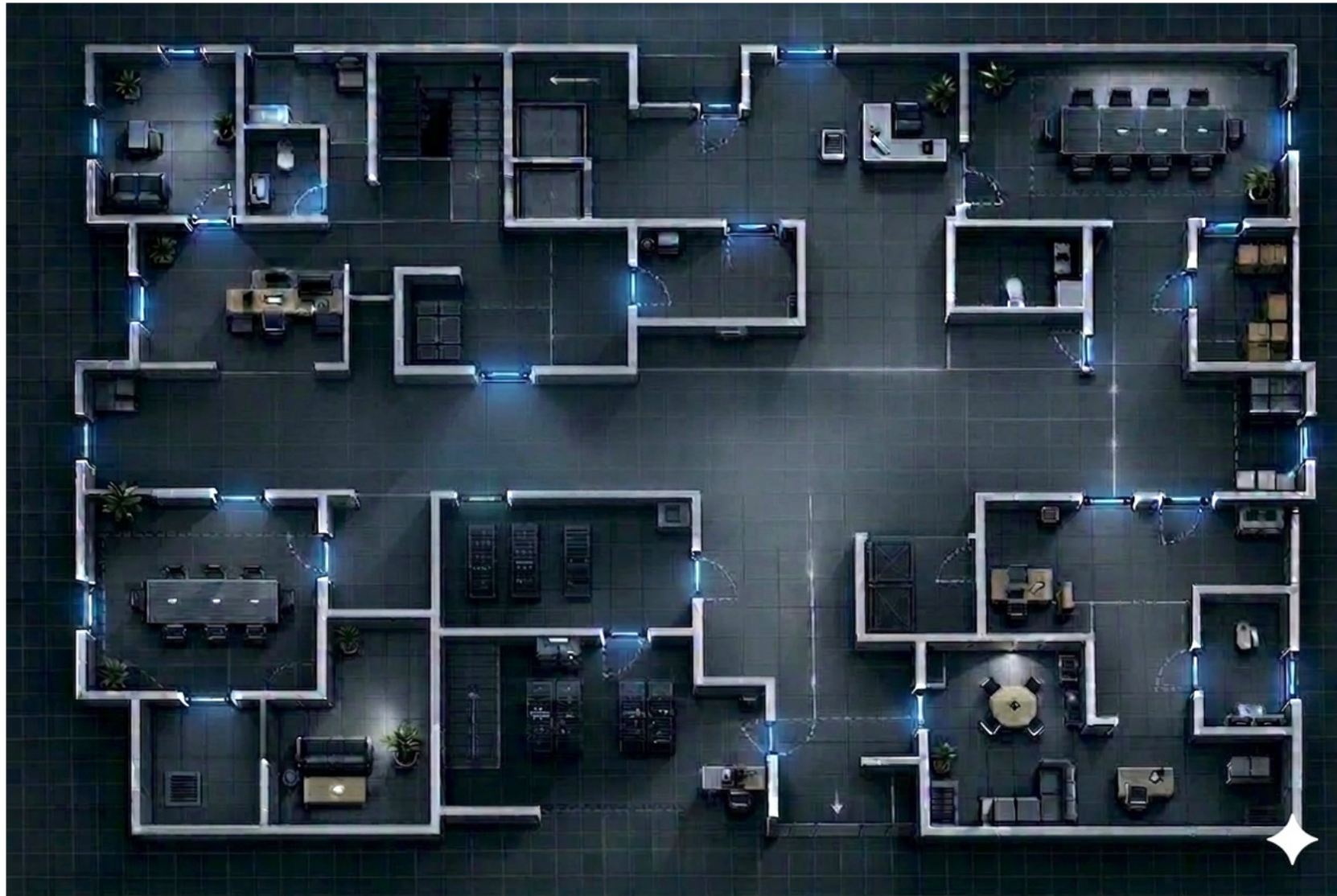
The Rise and Fall of Stack Overflow



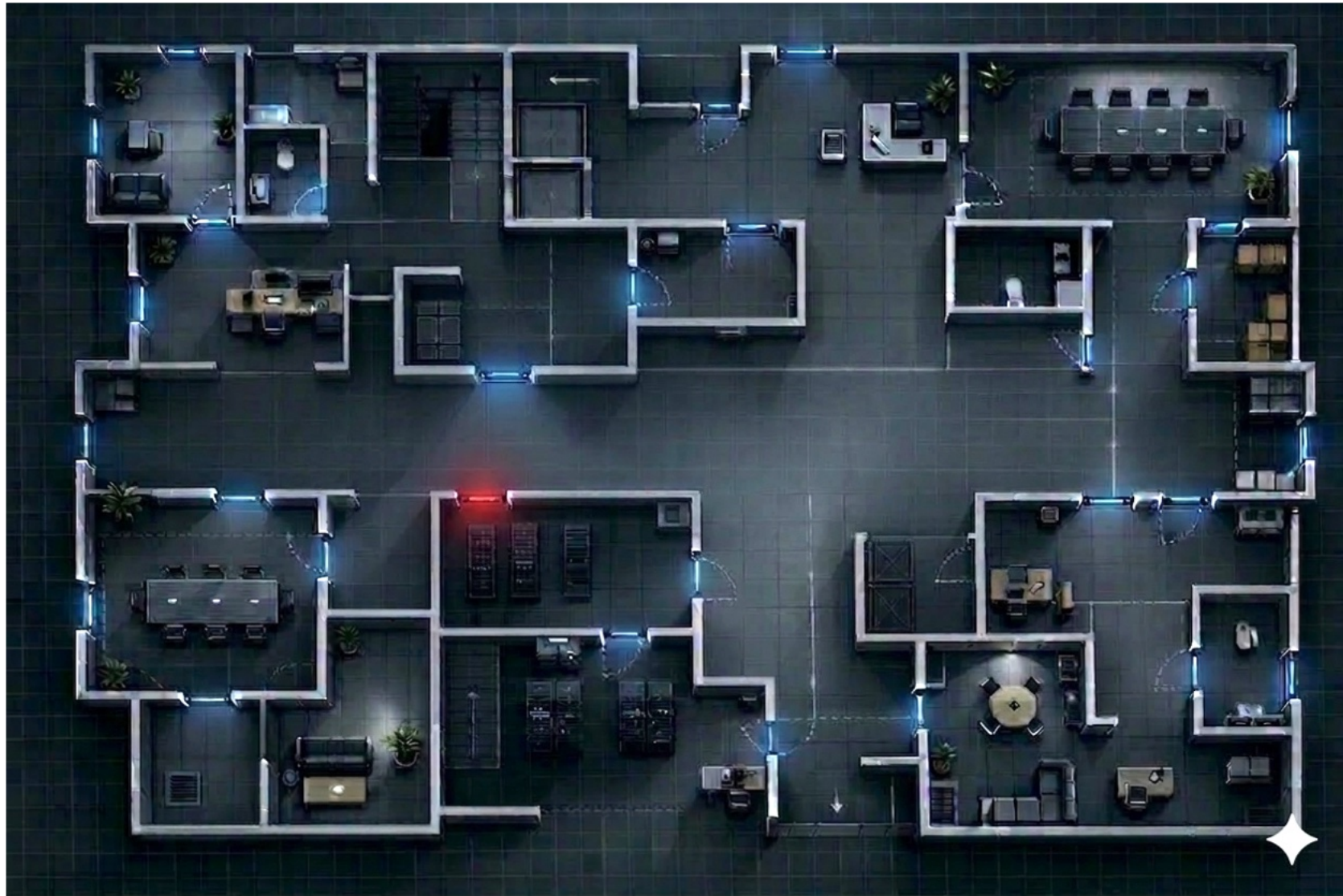
Hur påverkas cyberdomänen av AI?



Potentiella ingångar till nätverket



Vissa dörrar är olåsta



Vem hittar den öppna dörren först?



Snart får angriparna stor hjälp av AI



Försvararna kan också hitta sårbarheter med AI



Många sårbarheter hämmar försvararna



Prioritering avgörande för försvar



När AI också avhjälper av sårbarheter



Lugnet före stormen före lugnet

