



AI and Infrastructure

Prof. Fredrik Heintz

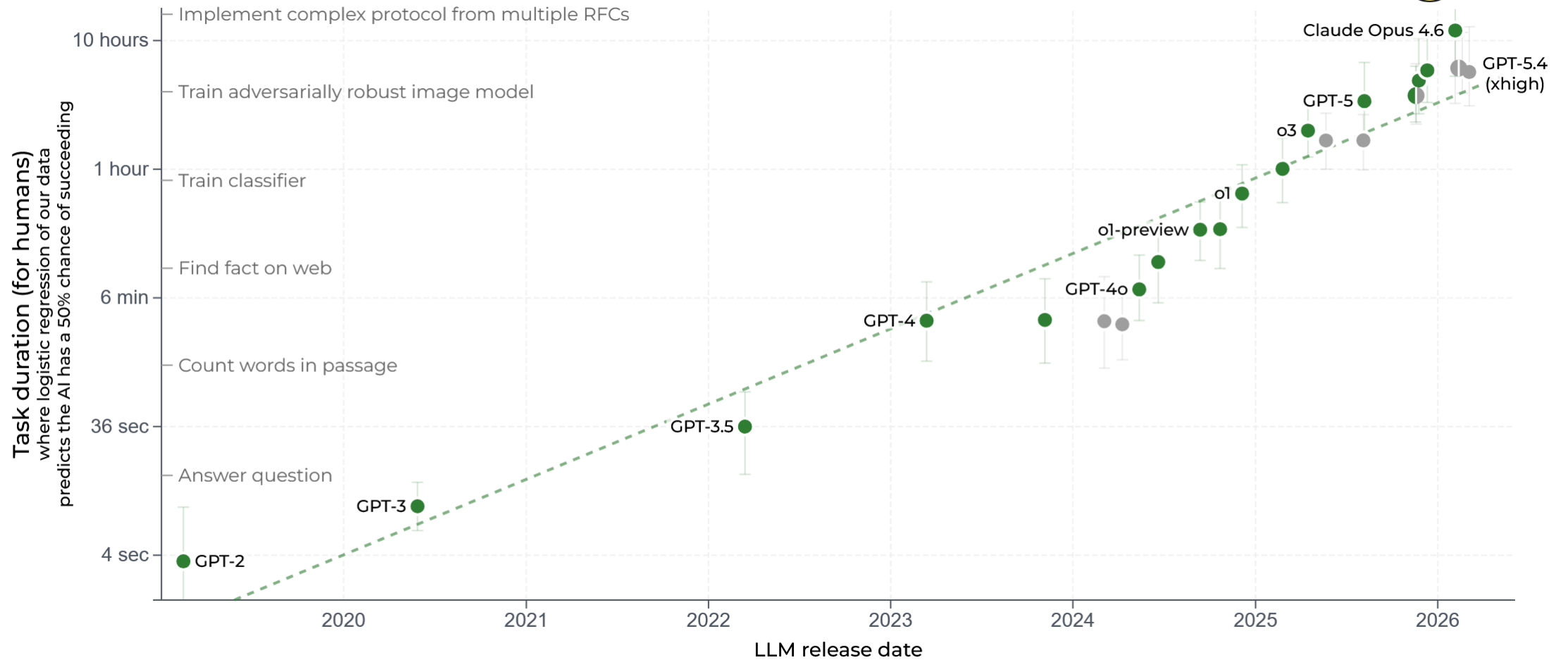
Institutionen för Datavetenskap, Linköpings universitet

Director AI4x Excellens Center

Co-Director WASP

fredrik.heintz@liu.se @FredrikHeintz

Time horizon of software tasks different LLMs can complete 50% of the time



Time Horizon 1.1 (Current) ▾

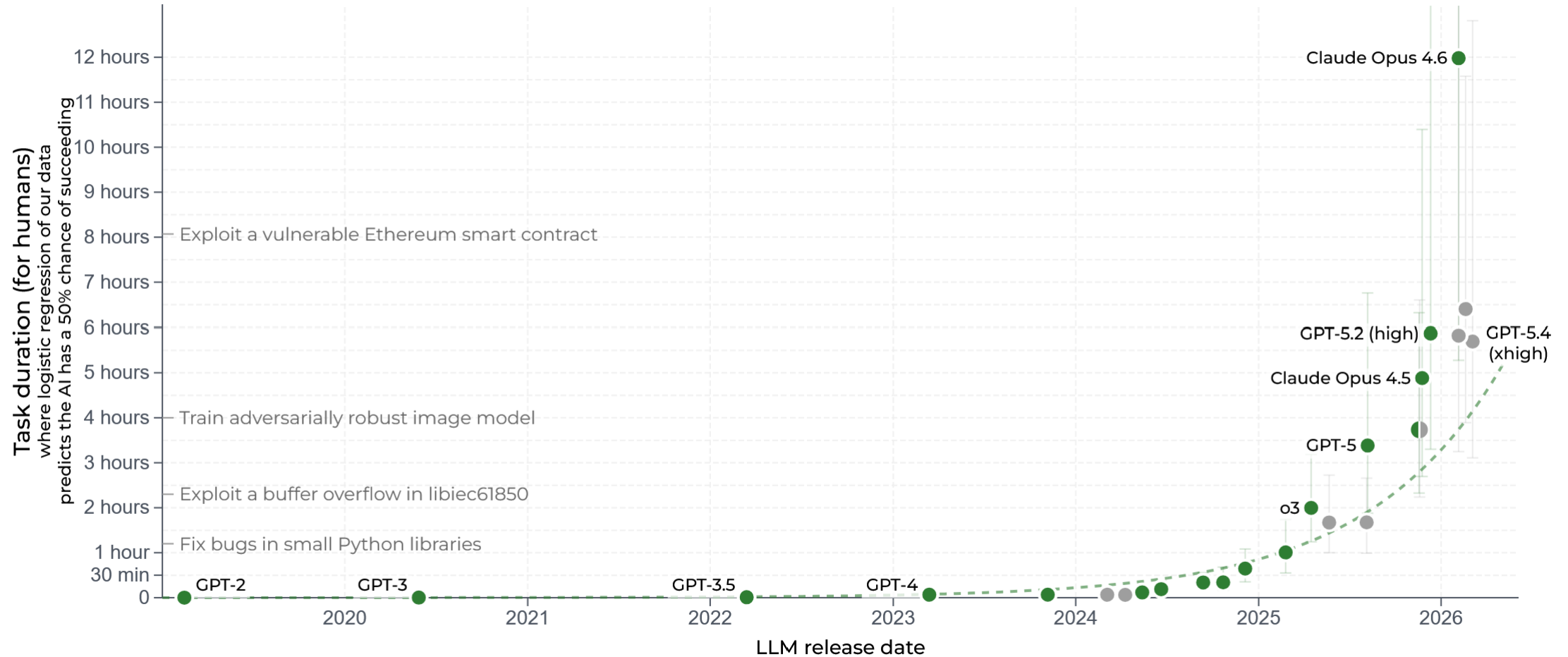
Log Scale Linear Scale

50% Success

80% Success



Time horizon of software tasks different LLMs can complete 50% of the time



Time Horizon 1.1 (Current) ▾

Log Scale

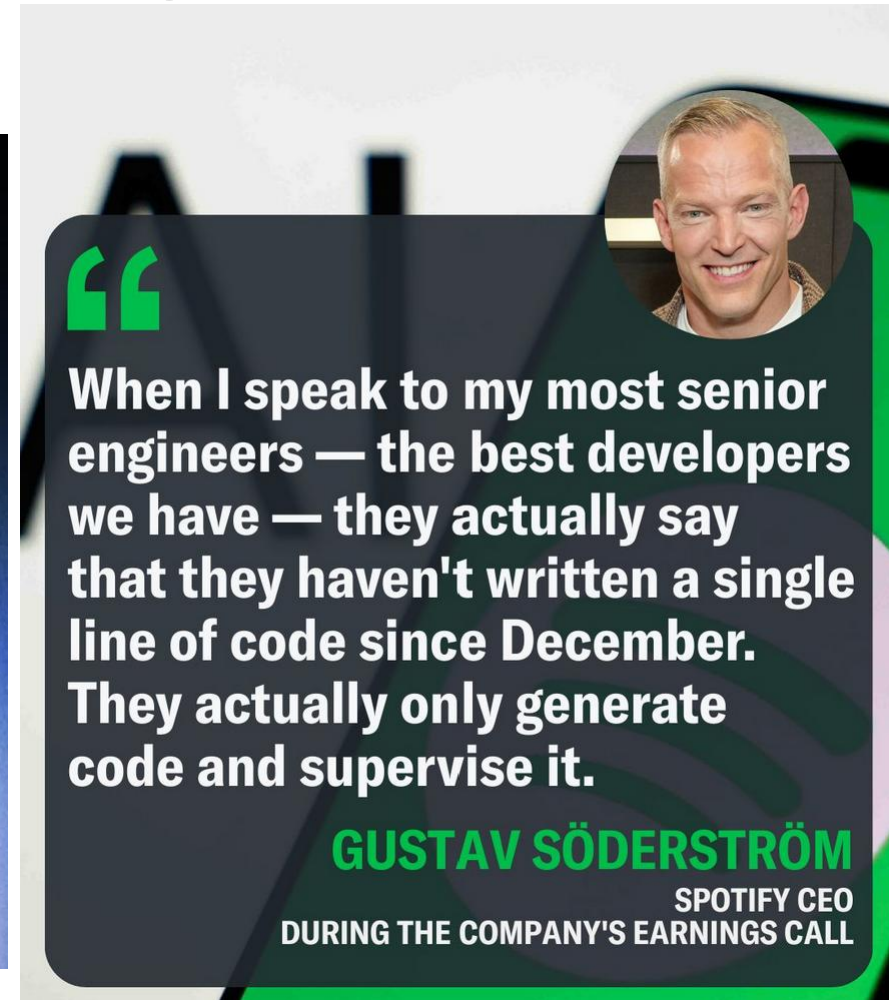
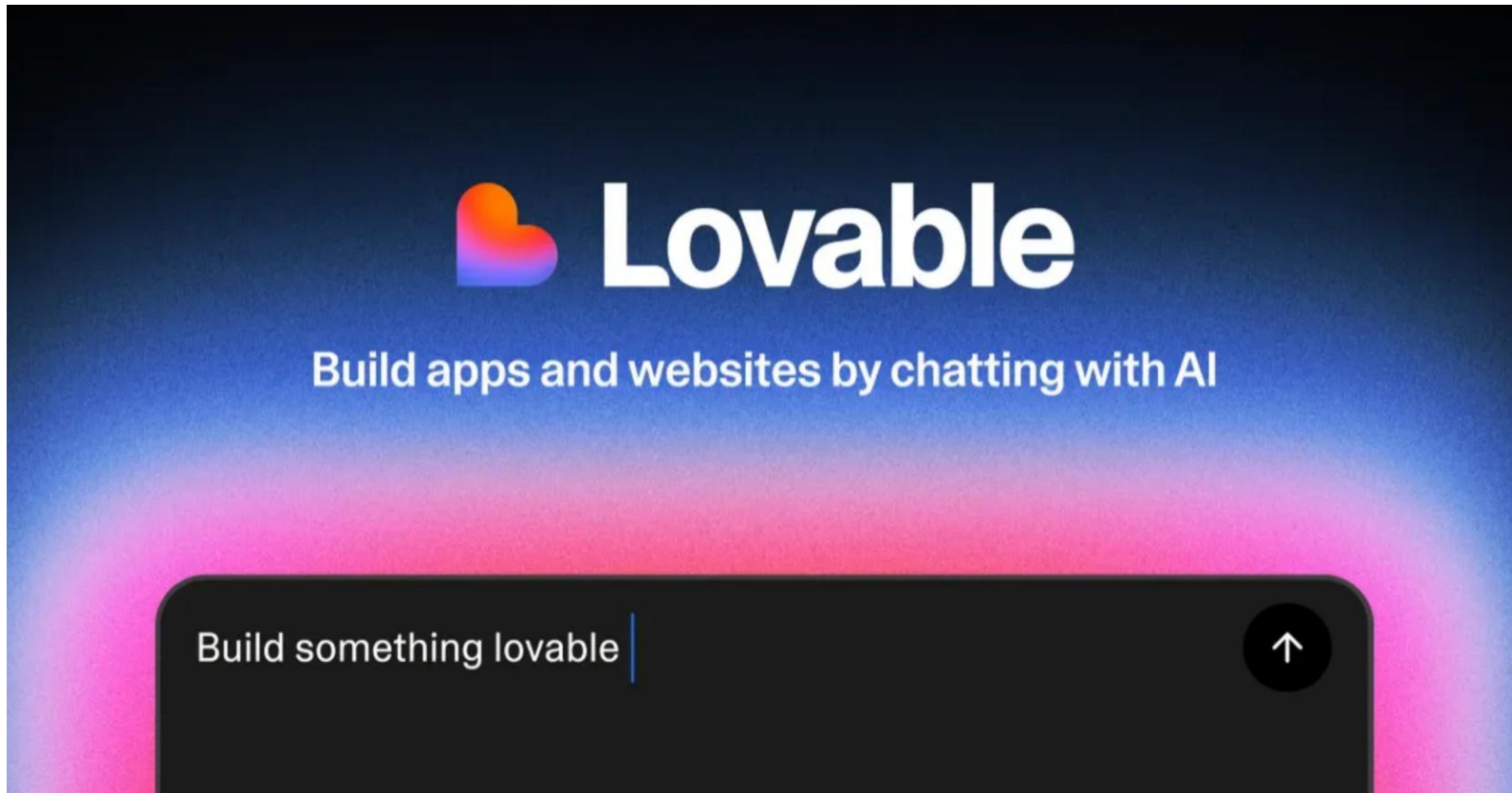
Linear Scale

50% Success

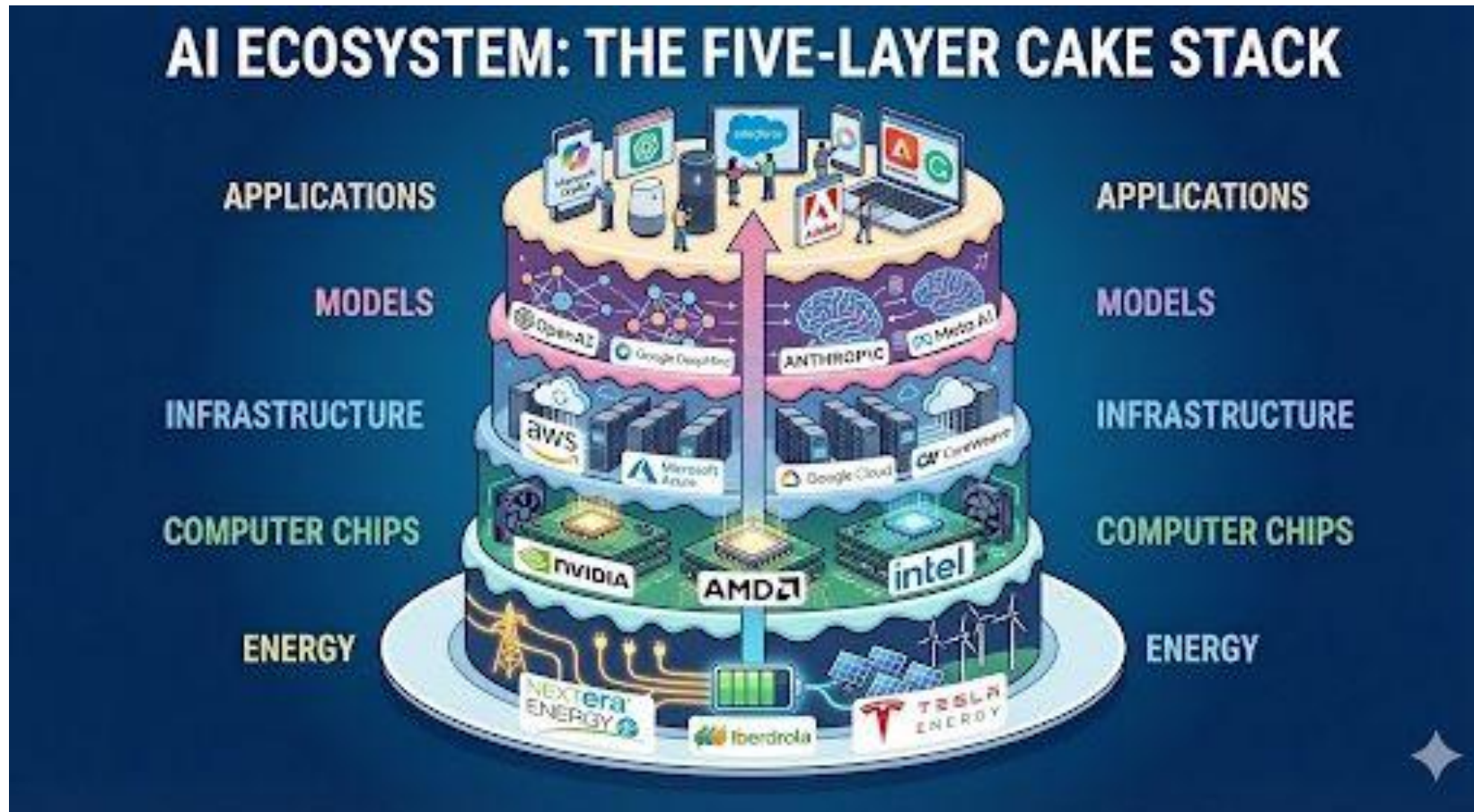
80% Success



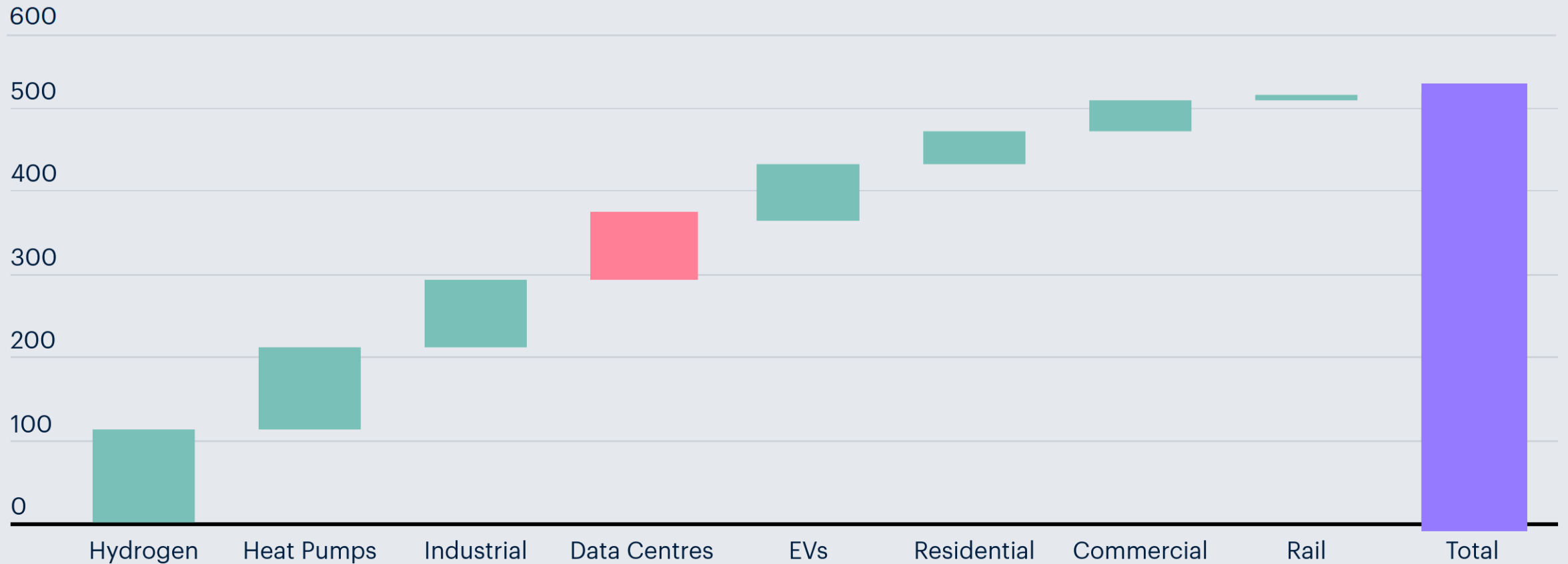
Coding Assistants and Vibe Coding



The Five Layer AI Stack (according to J. Huang)

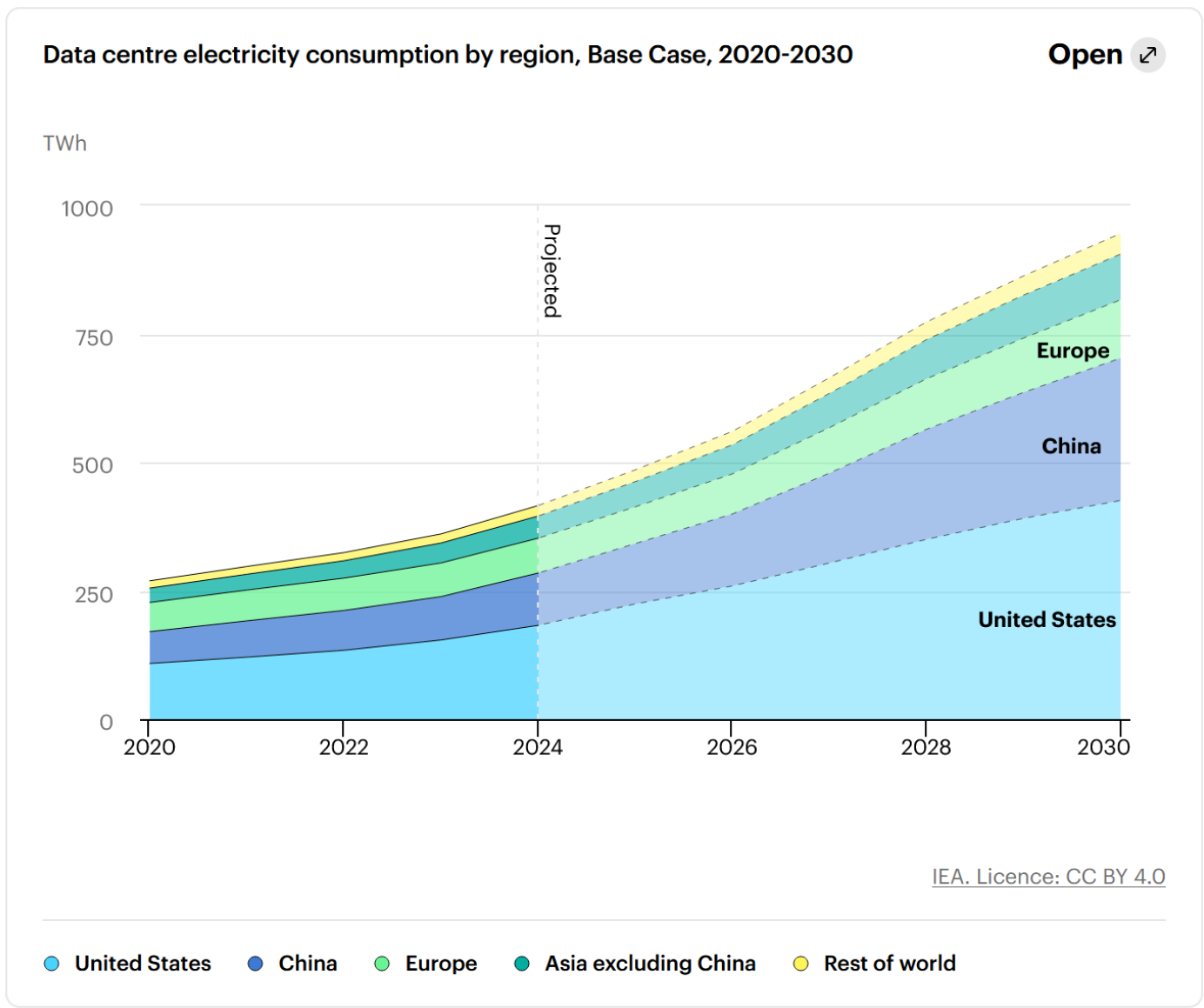


Change in European power demand by sector 2030 vs 2024, TWh



Source: ICIS

Note: Hydrogen demand is based on EU RFNBO targets, which are likely overly optimistic



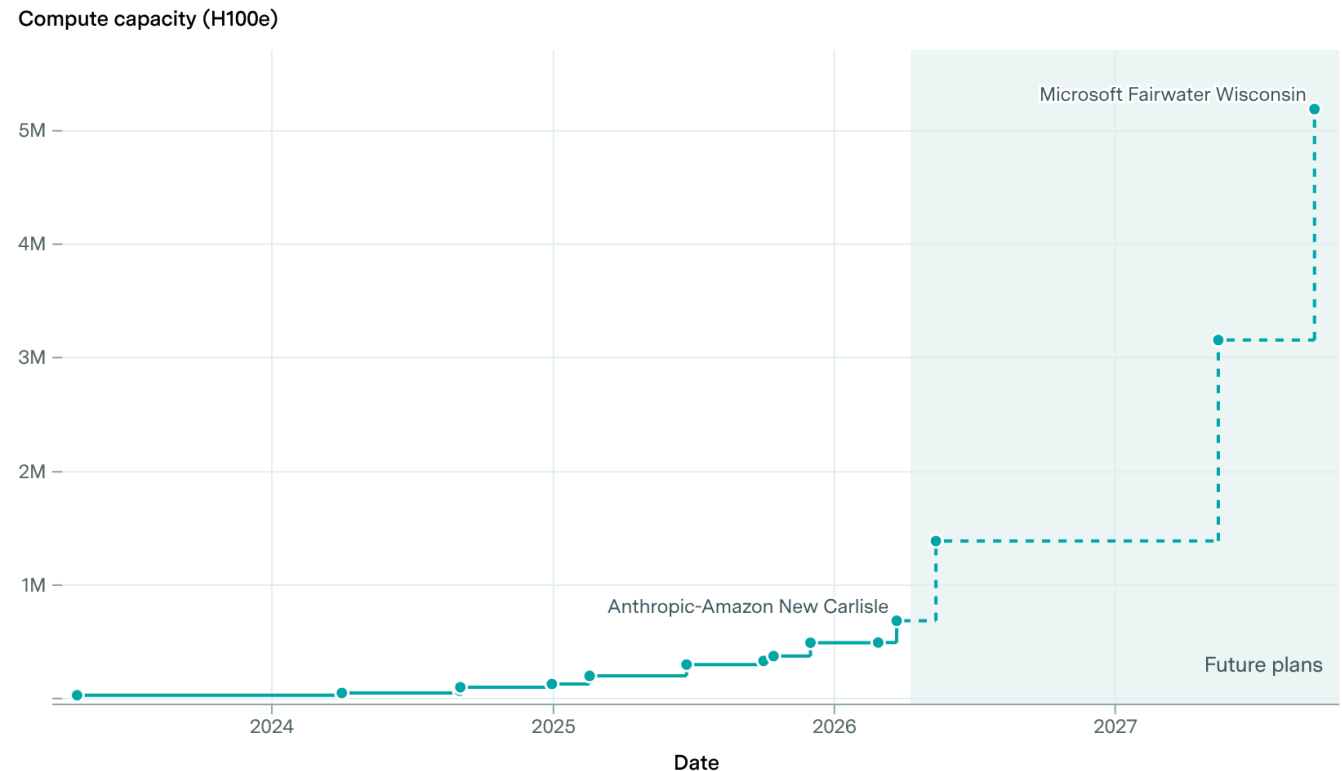
Compute Build-Out

● Data Centers

The largest known AI data center has computing power equivalent to 700,000 NVIDIA H100 chips.

Anthropic-Amazon New Carlisle has an estimated 1.1 GW of power capacity and \$35 billion in capital costs, making it the largest known AI data center. Microsoft's Fairwater Wisconsin data center is expected to be nearly 8x more powerful, at 5.2M H100-equivalents by September 2027.

[Learn more →](#)



EPOCH AI | CC-BY

[Download graph](#)

Training vs Inference

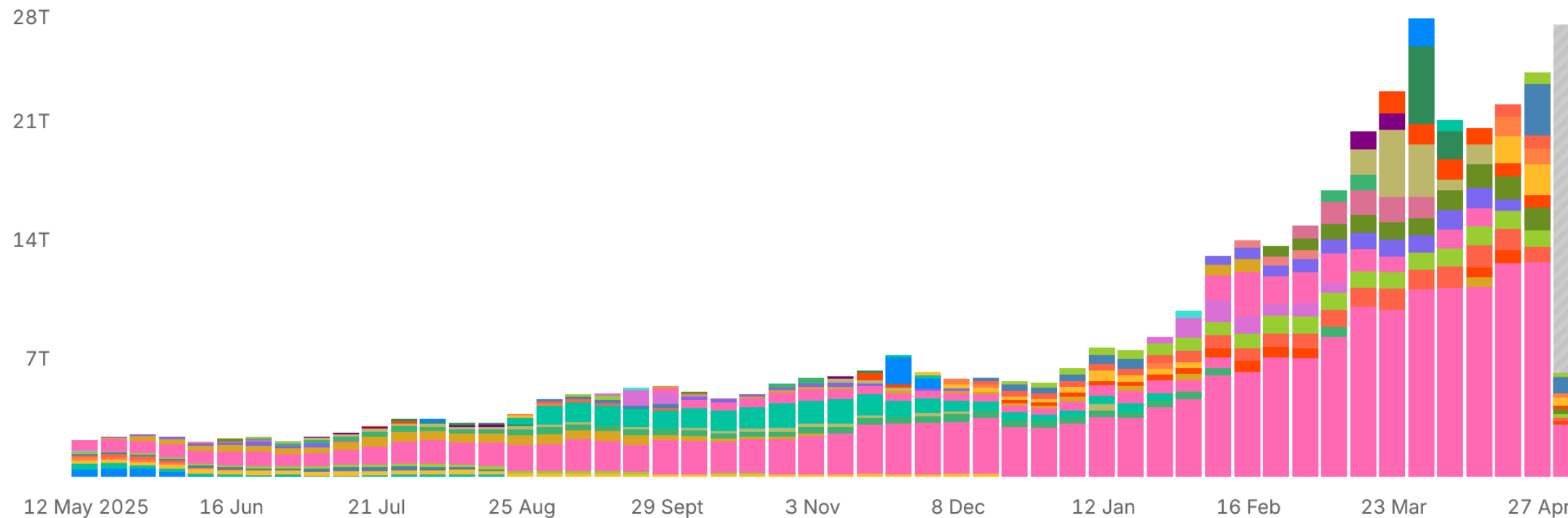
- **2023:** Inference accounted for roughly **33%** of AI compute.
- **2025:** Inference reached **50%** of total AI compute.
- **2026 (Projected):** Inference will reach **65% to 70%** of all AI compute.
- **2029-2030 (Projected):** Inference is expected to account for **70-90%** of total compute and infrastructure spend

IF YOUR WORKLOAD IS...	OPTIMIZE FOR	HARDWARE CHOICE	WHY
Training large models	Throughput	H100/H200, multi-node	Raw compute power matters
Production inference	Latency	B200/B300, specialized	User experience, cost per token
Variable inference load	Autoscaling	Cloud GPU instances	Match capacity to demand
Latency-critical inference	Edge deployment	Smaller GPUs distributed	Reduce network round-trip
Cost-sensitive inference	Efficiency	TPU, Trainium, AMD	30-40% savings possible

Token Usage

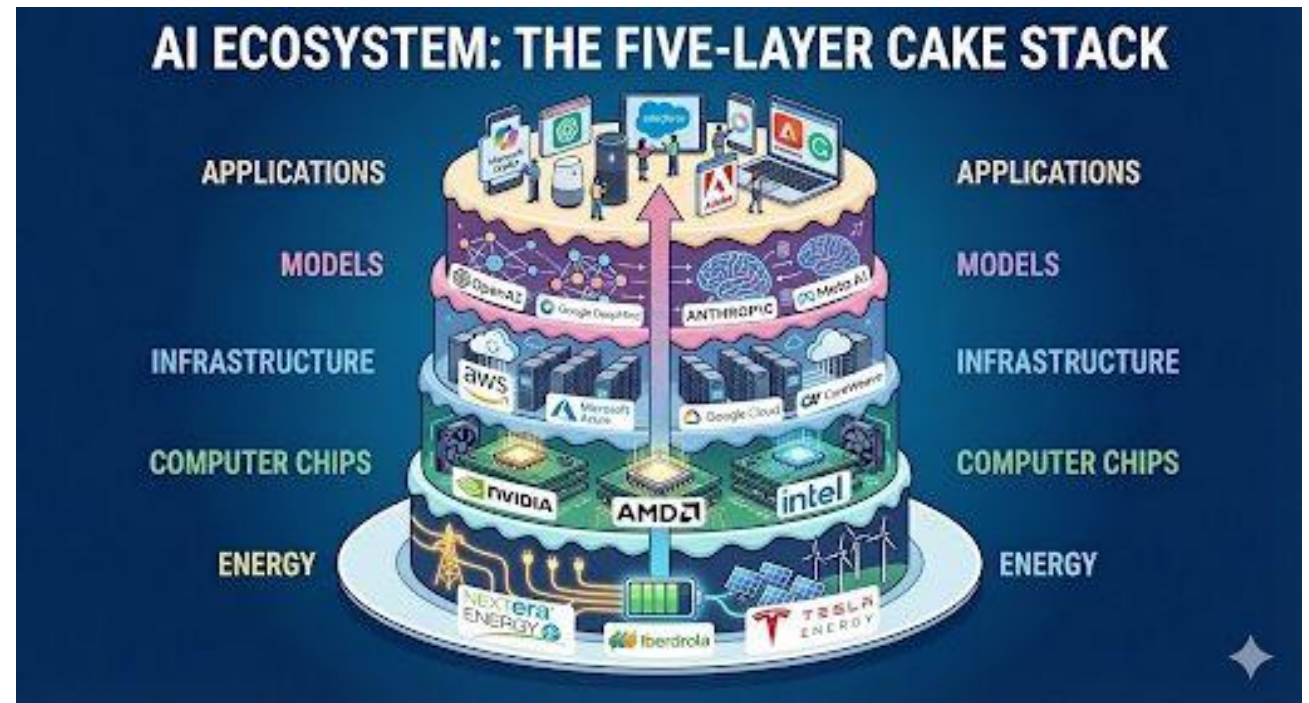
Top Models

Weekly usage of models across OpenRouter



Alternatives to the Current Paradigm

- Small models
- Distributed / federated models
- Neurosymbolic AI





Claude Mythos and Cyber Security, What the Leak Actually Tells Defenders



RESIST CRAI Center

Center Management (Director/PI, Manager, co-PIs)

Board

SAC

Research Program

Trustworthy and Verifiable AI
Runtime Security Assurance
Robust and Secure AI-Supported Development
Resilient Distributed and Agentic AI

Validation Program

AI-driven critical Infrastructure
AI-empowered software services
AI-automated robots and vehicles

Center Activities

Strategic research agenda
Education and training
Policy and strategy advise
Internationalization
Dissemination

Academia



LUNDS UNIVERSITET

Industry

COMBITECH



ERICSSON

SECTRA

SIEMENS
ENERGY



SAAB

Society



Internat.

IMPERIAL
Inria

The Alan Turing
Institute



ACTION AI INSTITUTE

Where are the Bottlenecks?

- Energy – Compute – Tokens
- Domain/Application Optimised Models
- Regulation, legal certainty
- Competence

