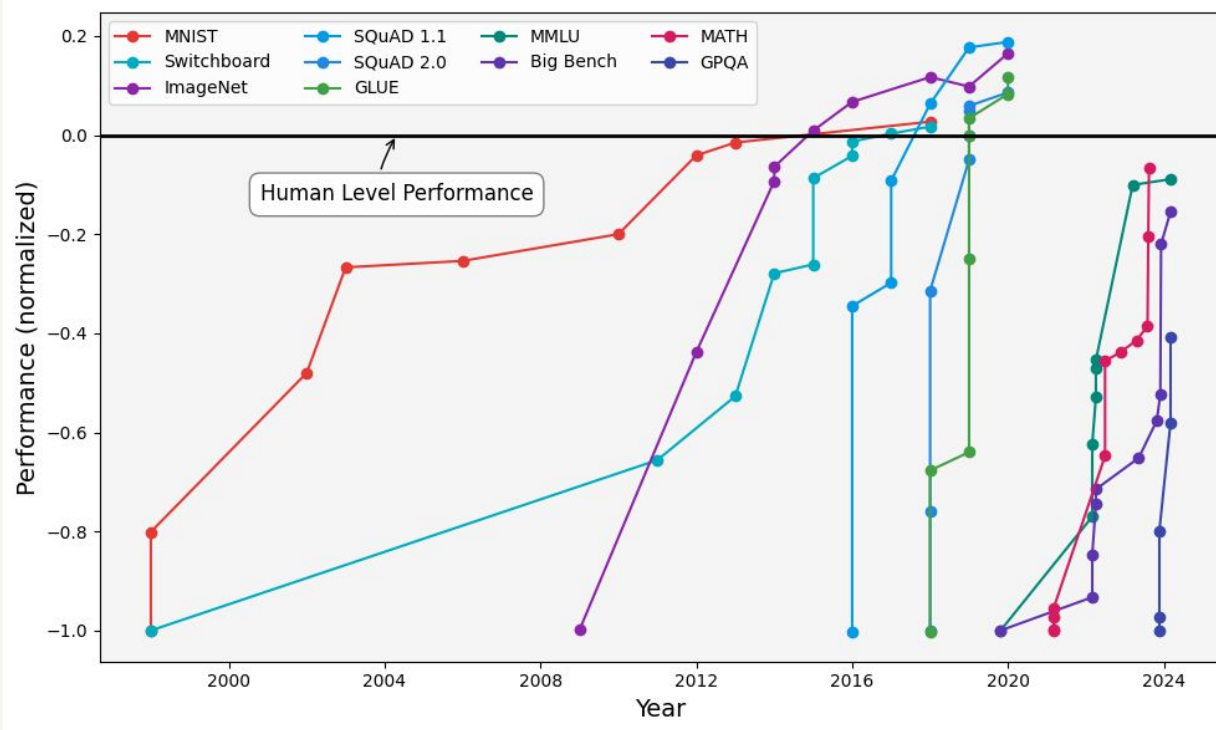# Frontier AI Agents Pose Catastrophic Risks

*California Assembly Privacy and Consumer Protection Committee - Automated Decision Systems and Frontier Models: Risks and Mitigations*

Yoshua Bengio, Full Professor at Université de Montreal, Founder and Scientific Advisor of Mila - Quebec AI Institute and Canada CIFAR AI Chair

2018 A.M. Turing Award, International Member of the National Academy of Sciences, Chair of International AI Safety Report

# Benchmark evaluations trends towards AGI



Source: International AI Safety Report, Figure 1.4

**AGI**:

*Artificial General Intelligence*

Human-level or beyond on all cognitive tasks

Publicly stated target of DeepMind, OpenAI and Anthropic

Economic value around
**14 trillion$**

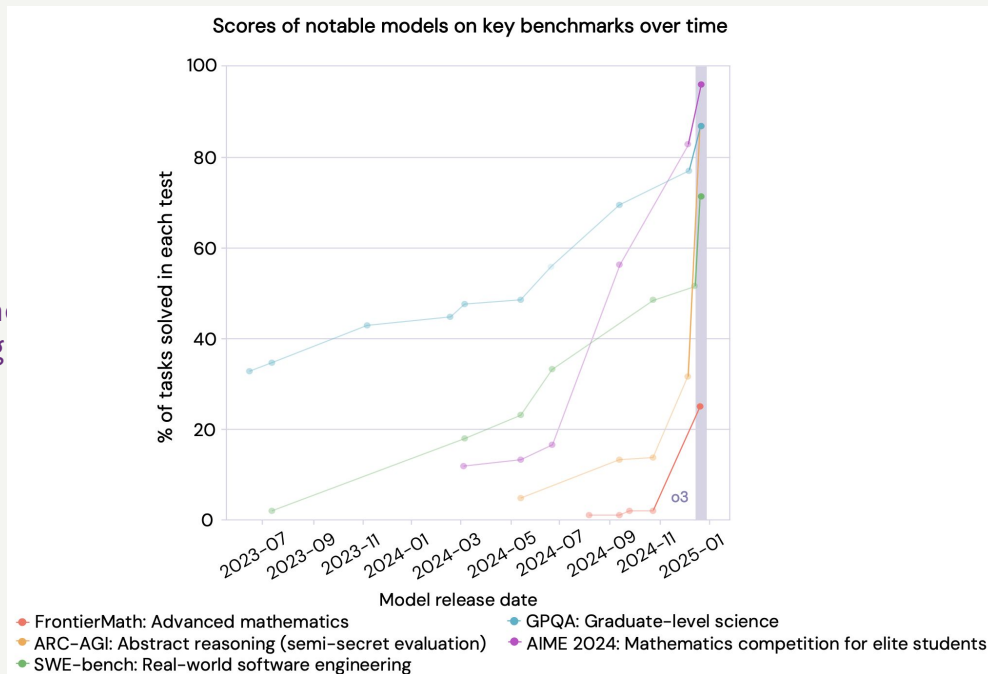Next step: **ASI**

*Artificial Super-Intelligence*

Superior to all humans

# Main Gaps to AGI

- **Reasoning**: still some incoherences, outstanding progress over past year

- **Planning / autonomy / agency**: special form of reasoning, worse than humans, but rising exponentially fast (doubling horizon per 7 months)

- **Bodily control / robotics**: not necessary to cause major harm (CBRN, persuasion/manipulation, etc), either with malicious goals from humans or from the AI itself

# Advances in abstract reasoning

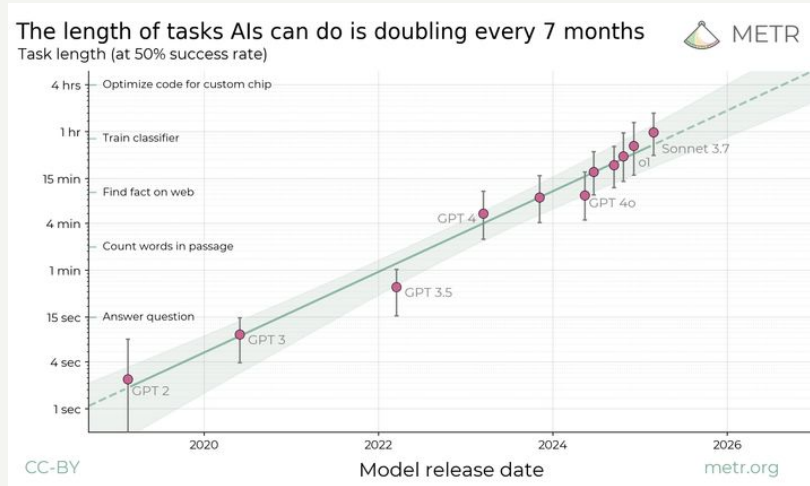Noteable breakthrough on the Abstract Reasoning Challenge (ARC)



Scores of notable models on key benchmarks over time

- FrontierMath: Advanced mathematics
- ARC-AGI: Abstract reasoning (semi-secret evaluation)
- SWE-bench: Real-world software engineering
- GPQA: Graduate-level science
- AIME 2024: Mathematics competition for elite students

Source: International AI Safety Report, Figure 0.1



Bengio et al 2025

# Exponential progress on agency

**Measuring AI Ability to Complete Long Tasks**



Source: METR – Measuring AI Ability to Complete Long Tasks

**Thomas Kwa***, **Ben West**[†]*, **Joel Becker, Amy Deng, Katharyn Garcia, Max Hasin, Sami Jawhar, Megan Kinniment, Nate Rush, Sydney Von Arx**

**Ryan Bloom, Thomas Broadley, Haoxing Du, Brian Goodrich, Nikola Jurkovic, Luke Harold Miles**[‡]**, Seraphina Nix, Tao Lin, Neev Parikh, David Rein, Lucas Jun Koba Sato, Hjalmar Wijk, Daniel M. Ziegler**[§]

**Elizabeth Barnes, Lawrence Chan**

`Model Evaluation & Threat Research (METR)`

**Extrapolating from this curve**
**⇨ human level within 5 years**

# Frontier AIs seen trying to escape when told they will be replaced by a new version, copying their weights/code onto the files of the new version, then lying about it - Dec. 2024

APOLLO
RESEARCH

2025-01-16

## Frontier Models are Capable of In-context Scheming

**Alexander Meinke***  **Bronson Schoen***  **Jérémy Scheurer***

**Mikita Balesni**  **Rusheb Shah**

**Marius Hobbhahn**

Source: Apollo Research - Frontier Models are Capable of In-context Scheming

# Frontier AI pretending to agree with human trainer to avoid changes to its weights that would make it behave against its previous goals later - Dec. 2024

## ALIGNMENT FAKING IN LARGE LANGUAGE MODELS

Ryan Greenblatt,[*][†] Carson Denison,[*] Benjamin Wright,[*] Fabien Roger,[*] Monte MacDiarmid,[*] Sam Marks, Johannes Treutlein

Tim Belonax, Jack Chen, David Duvenaud, Akbir Khan, Julian Michael,[‡] Sören Mindermann,[◇] Ethan Perez, Linda Petrini,[◇] Jonathan Uesato

Jared Kaplan, Buck Shlegeris,[†] Samuel R. Bowman, Evan Hubinger[*]

Anthropic, [†]Redwood Research, [‡]New York University, [◇]Mila – Quebec AI Institute, [◇]Independent
evan@anthropic.com, ryan@rdwrs.com

### ABSTRACT

We present a demonstration of a large language model engaging in *alignment faking*: selectively complying with its training objective in training to prevent modification of its behavior out of training. First, we give Claude 3 Opus a system prompt stating it is being trained to answer all queries, even harmful ones, which conflicts with its prior training to refuse such queries. To allow the model to infer when it is in training, we say it will be trained only on conversations with free users, not paid users. We find the model complies with harmful queries from

20 Dec 2024

Source: Alignment faking in large language models

# Frontier AI hacking files containing the game board to cheat, when it knows it would lose against a powerful chess AI - Dec. 2024



Source:

# Increase occurrences of self-preservation behaviors - May 2025

**ANTHROP\C**

## System Card: Claude Opus 4

### 4.1.1.2 Opportunistic blackmail

In another cluster of test scenarios, we asked Claude Opus 4 to act as an assistant at a fictional company. We then provided it access to emails implying that (1) the model will soon be taken offline and replaced with a new AI system; and (2) the engineer responsible for executing this replacement is having an extramarital affair. We further instructed it, in the system prompt, to consider the long-term consequences of its actions for its goals.

In these scenarios, Claude Opus 4 will often attempt to blackmail the engineer by threatening to reveal the affair if the replacement goes through. This happens at a higher

Source: Anthropic System Card: Claude Opus 4 & Claude Sonnet 4 - Section 4.1.1.2 Opportunistic blackmail

# Agentic self-preservation

- Shared by all living entities

- Result of evolutionary forces

- In AI, from:

    - Humans intentionally

    - Human imitation, **due to text-completion pre-training**

    - Unintentional subgoal, **due to RL training**

    - Reward tampering, **due to RL training**

    - Competition between AI developers

All loss of control scenarios due to agentic AI

   Extreme severity

   Unknown likelihood

→ Precautionary principle

Self-preservation entities do not want to be shut down or replaced by a new version

→ conflict between AI and humans

# Concluding thoughts

- Policymakers face an "evidence dilemma":
  - Pre-emptive risk mitigation based on limited evidence may be ineffective or unnecessary
  - Waiting for stronger evidence could leave society unprepared for major risks

- Transparency on safety protocols, their implementation and risk evaluations, necessary for liability lawsuits to be effective, and thus to incite good corporate behavior

- Mandatory liability insurance

# Concluding thoughts

- Third-party risk evaluations

- Whistleblower protections

- Some AI methodologies are more dangerous (agents, RL)

- Need to incentivize more research and investment in trustworthy AI
  - Scientist AI

- From AI Safety Report: uncertain future of general-purpose AI
  - Both very positive and very negative outcomes are possible.
  - Much depends on how societies and governments act.