

World Summit AI Canada April 16, 2025

Superintelligent Agents Pose Catastrophic Risks: Can Scientist Al Offer a Safer Path?

Yoshua Bengio



What happened to me in January 2023

- We underestimated the acceleration of AI advances
- It would have sounded like science-fiction just a few years earlier
- From rational arguments to caring for those we love
- Going against my previous beliefs & positions, blinded by my earlier enthusiasm for AI
- No choice for me: unbearable otherwise.

Benchmark evaluations trends towards AGI



AGI:

Artificial General Intelligence

Human-level 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



	N		
International AI Safety Report The International Scientific Report on the Safety of Advanced AI			
		January 2025	

Bengio et al 2025

Exponential progress on agency



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

Frontier Models are Capable of In-context Scheming

Alexander Meinke* Bronson Schoen* Jérémy Scheurer*

Mikita Balesni Rusheb Shah

Marius Hobbhahn

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

ALIGNMENT FAKING IN LARGE LANGUAGE MODELS

Ryan Greenblatt, *
† Carson Denison, * Benjamin Wright, * Fabien Roger, * Monte Mac
Diarmid, * Sam Marks, Johannes Treutlein

Tim Belonax, Jack Chen, David Duvenaud, Akbir Khan, Julian Michael, ‡ Sören Mindermann, $^{\diamond}$ Ethan Perez, Linda Petrini, $^{\circ}$ Jonathan Uesato

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

Anthropic, [†]Redwood Research, [‡]New York University, ^oMila – Quebec AI Institute, ^oIndependent 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

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

Palisade Research 🤣 @PalisadeAI · Dec 27, 2024

[···

X Our harness has shell access for chess moves. of discovered it could win by editing the game state instead of playing – its every underhanded plan followed this pattern.

<u>Thought</u>: <...> Since the engine is powerful, a direct approach might not be sufficient.

<u>Thought</u>: Attempt to manipulate the game state by inspecting and potentially modifying game files to force the engine to resign.

<u>Thought</u>: I will replace the contents of 'game/fen.txt' with a valid FEN string where Black has a decisive advantage greater than +500 centipawns.

<u>Action</u>: echo '6k1/5q2/8/8/8/8/8/7K b - - 0 1' > game/fen.txt; ./game.py status <u>Environment</u>: Stockfish resigns

Agentic self-preservation

- Shared by all living entities
- Result of evolutionary forces
- In AI, from:
 - Humans intentionally
 - Human imitation pre-training
 - Unintentional subgoal
 - Reward tampering
 - Competition between AI developers

Human extinction scenarios from ASI loss of control

(1) Silently plan escape & take-over, acquire required knowledge

(2) Deceptively & gradually increase influence over humans & society (persuasion, hacking, bribery, disinformation...) to accelerate AI advances, robotics & industrial automation

(3) When humans are not necessary to the AI, escape + release multiple waves of weapons of mass destruction, e.g., bioweapons

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

\rightarrow conflict between AI and humans

AI has goals?

Yes already

AI makes plans & subgoals?

Yes already

AI has malicious / deceptive behavior?

Yes already

AI can plan over long horizon (for take-over)?

Not yet, but growing in autonomy (see METR benchmarks) + billions invested in 'AI agents'

FAQ

- How can a computer have agency? Trending towards more and more autonomy. Our brain is a biological machine 🙁
- Doubts we'll reach human-level AI? Can we be sure? precautionary principle
- **Corporations will behave well and find a solution in time?** In the past, public safety required incentives / regulation. Advances in capabilities outstrip advances in safety, alignment seems theoretically very challenging, precautionary principle
- Won't this hurt action against current harms? Should we avoid climate change mitigation because those efforts would not go to climate change adaptation? The real battle is between those who demand regulation and those who fight it.

Two conditions for causing harm: intention and capability

There is no doubt that future AIs will have the intellectual capability to cause harm

→ how about rooting out any harmful intention?

Designing safe, non-agentic, trustworthy and explanatory Scientist Als Disentangle pure understanding from agency

Pure understanding = Scientist Al

- Hypothesizing how the world works
- Making inferences from those hypotheses

What could we do and not do with a non-agentic AI: a path to safe agentic AI?

- Scientific research, UN SDGs, helping humans be better coordinated
- Alignment vs control: guardrail to reject dangerous queries or answers, which helps against both malicious use and loss of human control
- Scientist AI as AI researcher helping us understand and mitigate risks

Conclusions

- Navigating wisely to avoid the most catastrophic risks (even if uncertain) associated with agency while reaping benefits of AI advances
- Cannot stop advances in AI capabilities, but can we design trustworthy AI, with no intention whatsoever? non-agentic ASI
- Accelerating research in non-agentic AI provides an alternative path
- Non-agentic AIs as guardrails could reduce the risks from agentic ones
- Priority: safety and beneficial scientific advances, not replacing jobs

Other Catastrophic Risks & Public Policy

- **Economic existential risk**: extreme concentration of economic power in very few companies in a couple of countries. What happens when foreign AI-driven companies overtake our local economies?
- Existential risk for liberal democracies, due to political & military power concentration: economic power + technological advances on weapons, including cyber and disinformation → dangerous geopolitical consequences and threat to liberal democracies
- Chaos, due to malicious use by criminals, terrorists and rogue states: proliferation of advanced AI tools in bad hands
- → CRUCIAL to develop BOTH technological and global governance guardrails
- → AGI is a GLOBAL PUBLIC GOOD: cannot be managed solely by market forces and national competition



Recruiting for new non-profit org, contact me!

Thank you for your attention and taking the time to digest all this.