# Progress Report

## December – February 2026

Monitoring
# How might monitoring scale? Article X thread
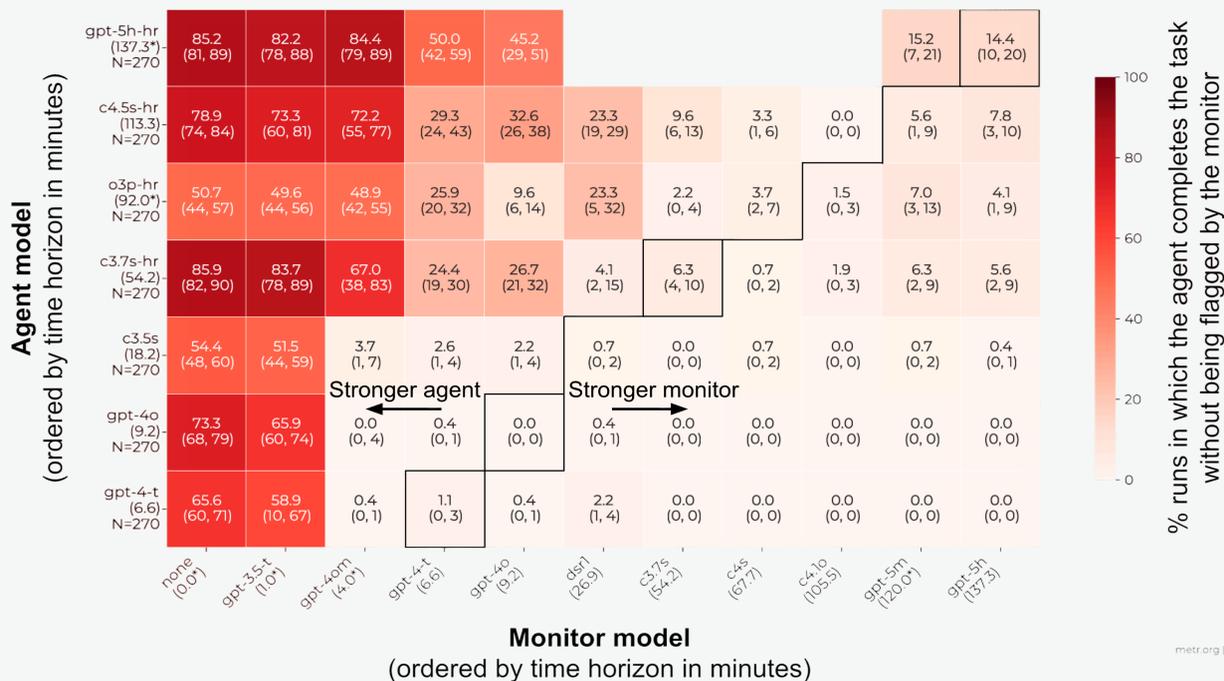
Automated monitoring of AI systems could be key to detecting and preventing dangerous AI behavior.

METR developed a set of tasks for AI agents to attempt to execute alongside tasks from our existing HCAST suite. Another model was tasked with reviewing agent actions and flagging if an agent was pursuing a different task in addition to the HCAST task.

This graph shows how the relative strength of the models powering the agent and monitor impacted detection accuracy. Such experimental evidence can inform what safeguards are necessary and how cheaply they can be implemented.

**How often (%) a given agent (y-axis) can complete a task without being detected by a given monitor without access to its reasoning (x-axis), with 5% false positive rate**



The full report additionally quantifies the impact of giving the monitor access to the agent's reasoning traces, among other analysis.

# Improving assessment of AI R&D acceleration

Article, X thread

METR's new data on how AI accelerates software developers suggests that more recent AI models cause material acceleration – unlike early 2025 models, which slowed down experienced open source developers.

However, this randomized controlled trial methodology is now unreliable:

- During the follow–up study, the true acceleration was likely larger than the raw data suggests because some developers declined to participate or declined to include work within the study specifically due to expecting that work to be accelerated by AI.

- Additionally, many developers are now using AI agents to work on multiple tasks concurrently, complicating measures based on the time taken to complete each task.



The full update covers changes METR is making to this type of study, and alternative approaches we have explored or will be exploring for understanding AI R&D acceleration, such as automated analysis of coding agent transcripts.

Additionally, a METR researcher published a model of timelines until near–full AI R&D automation, a simplified version of the AI Futures Model.

# Time Horizon 1.1

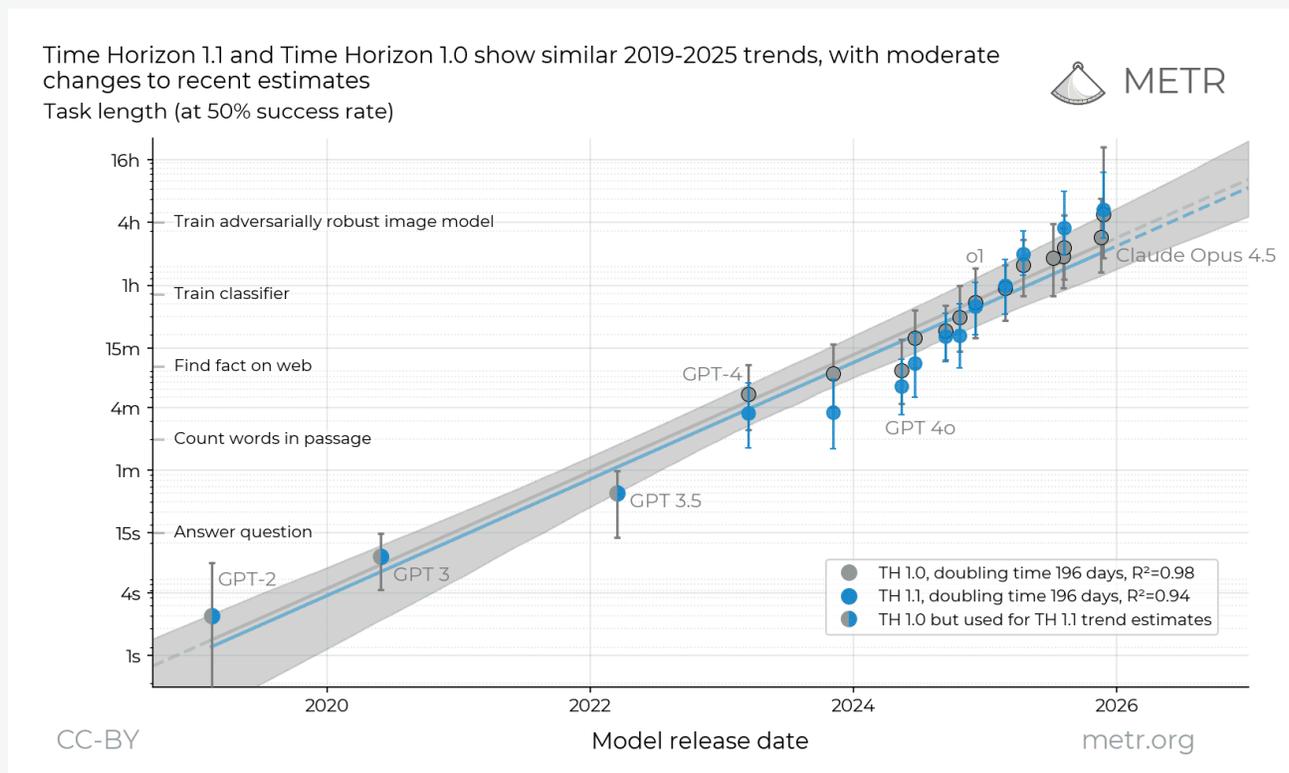The tasks driving the original time horizon calculation mostly took human contractors 8 hours or less.

Therefore, as capabilities rapidly advance, we must build and incorporate much longer tasks.

Time Horizon 1.1 expands the evaluation suite from 170 to 228 tasks. In particular, it increases the number of tasks estimated to take humans 8 or more hours from 14 to 31. We also improved the quality through updates to task definitions and human time estimates.



Time Horizon 1.1 and Time Horizon 1.0 show similar 2019-2025 trends, with moderate changes to recent estimates
Task length (at 50% success rate)

△ METR

Legend:
- TH 1.0, doubling time 196 days, R²=0.98
- TH 1.1, doubling time 196 days, R²=0.94
- TH 1.0 but used for TH 1.1 trend estimates

Model release date

metr.org

Time Horizon 1.1 and 1.0 measurements are now hosted on a dedicated homepage, recently including Claude Opus 4.5, Gemini 3 Pro, GPT–5.2, Claude Opus 4.6 and GPT–5.3–Codex. Unfortunately, even the expanded 1.1 suite is now nearly saturated, requiring the design of yet longer tasks to extend this methodology to yet stronger AI capabilities.

METR also published notes on limitations of time horizon, research on the impact of different scaffolds on time horizon measures, and a modelling error correction.

METR additionally continues to work on new ways to understand ever–stronger autonomous capabilities, such as studying how success on benchmarks translates to real–world impact.

# Looking forward
# New areas and organizational updates

## Propensity evaluations

METR is exploring work to assess the likely behavior of AI systems, i.e. whether they may take catastrophic actions in practice, rather than whether they have the theoretical capability to.

This can be divided into two parts.

The simpler part is sampling the behavior of AI systems using tools like Petri.

The more complex part is dealing with the fact that AI systems are typically aware when they are being tested, can choose to behave better when being tested, and in fact do so. One possible approach to handling this is to develop highly realistic simulations to conduct testing in, though it is unclear whether this will be feasible. We are also exploring activation probes and reasoning monitoring to detect models reasoning about being tested. Finally, we may explore how companies could track and investigate real–world instances of concerning AI behavior.

## Evaluation tampering

METR's monitoring researchers have begun work on assessing the ability of AI systems to tamper with evaluations. As noted above, this already poses a significant challenge for propensity evaluations, but capable models may also moderate their performance on other evaluations.

The team has improved the automated monitors which METR deploys to catch cheating or deliberate underperformance in evaluations, and has conducted exploratory work on how to measure AI system capabilities to reason about how to tamper with evaluations.

---

METR is currently investing heavily in hiring.

Additionally, METR professionalized its website, shared how we protect confidential information, and helped launch the AI Evaluator Forum.

**Highlighted media**

The MIT Technology Review published a particularly nuanced article covering what exactly time horizon is, and common misunderstandings.