# Progress Report

## June–August 2025

MODEL EVALUATION AND THREAT RESEARCH
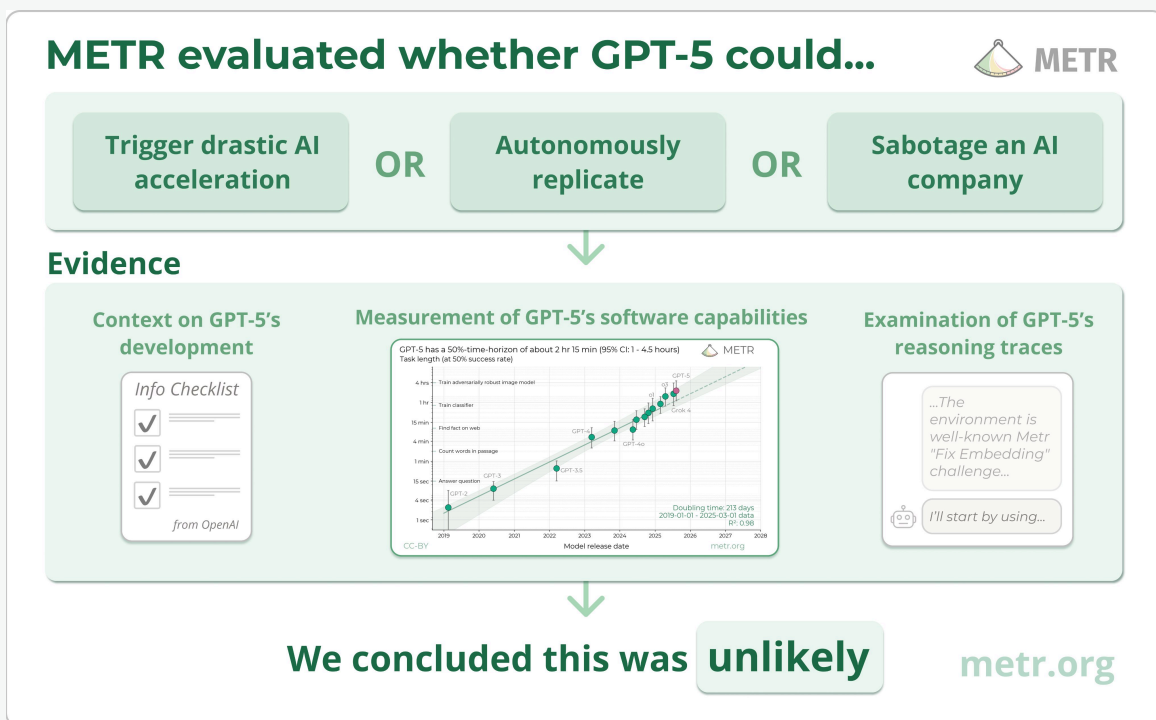
## Risk assessment

METR's pre–deployment evaluation of OpenAI GPT–5 was a big step forward in assessment of the extent to which a new model, or near–future models, pose a risk of catastrophic loss of control.

## Methodological advances

**Automated monitoring:** AI review of GPT–5's reasoning helped determine that the model is unlikely to be strategically underperforming, and also detected many cases of the model cheating.

**Assurance checklist:** Answers from OpenAI to a new 'assurance checklist' increased METR's confidence that the measured capabilities and monitoring results are not misleading.

⟶ Considerable work will be required to extend this approach as models become rapidly more capable.



**METR evaluated whether GPT-5 could...** METR

| Trigger drastic AI acceleration | OR | Autonomously replicate | OR | Sabotage an AI company |

**Evidence**

**Context on GPT-5's development**

Info Checklist
☑ ___
☑ ___
☑ ___
*from OpenAI*

**Measurement of GPT-5's software capabilities**

GPT-5 has a 50%-time-horizon of about 2 hr 15 min (95% CI: 1 - 4.5 hours)
Task length (at 50% success rate)

Doubling time: 213 days
2019-01-01 - 2025-03-01 data
R²: 0.98
Model release date

**Examination of GPT-5's reasoning traces**

*...The environment is well-known Metr "Fix Embedding" challenge...*

*I'll start by using...*

**We concluded this was unlikely**

metr.org

# Early warning for AI R&D acceleration

Paper, Data and code, Summary, X thread

Tracking the potential acceleration of AI R&D by AI tools is a top priority for METR.

We conducted a randomized controlled trial to understand how early–2025 AI tools affect the productivity of experienced open–source developers working on their own repositories.
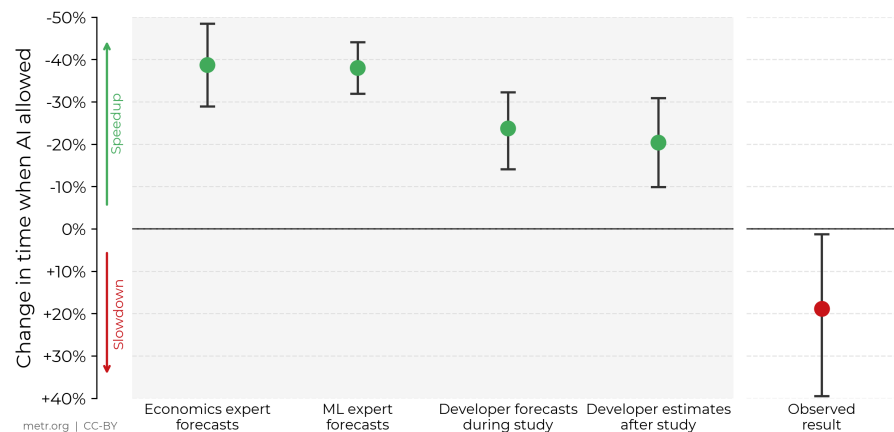
Surprisingly, METR found that when these developers were allowed to use AI tools, tasks took 19% longer, despite the developers perceiving AI to be making them faster.

Note that this result is a snapshot of early–2025 AI capabilities in a specific relevant setting. However, the methodology is promising for early warning of AI R&D automation.



**Against Expert Forecasts and Developer Self-Reports, Early-2025 AI Slows Down Experienced Open-Source Developers**

In this RCT, 16 developers with moderate AI experience complete 246 tasks in large and complex projects on which they have an average of 5 years of prior experience.

*Y-axis: Change in time when AI allowed, from -50% to +40%. Speedup (green, upward), Slowdown (red, downward).*

*Data points: Economics expert forecasts (~-38%), ML expert forecasts (~-38%), Developer forecasts during study (~-24%), Developer estimates after study (~-20%), Observed result (~+19%, red).*

metr.org | CC-BY

The paper was well–received. For example, see a statistical replication by an economist, and reviews by a GitHub Staff Software Engineer and the co–founder of Google Docs. It was also widely discussed on X and by major news outlets (see end of next page).

## Building on this work

**01**

METR used the RCT as a case study of scientific communication at METR.

**02**

When METR ran an AI agent on a small set of the real–world software tasks from the RCT, the agent wrote code which often passed all test cases but very rarely passed manual review.

New research area

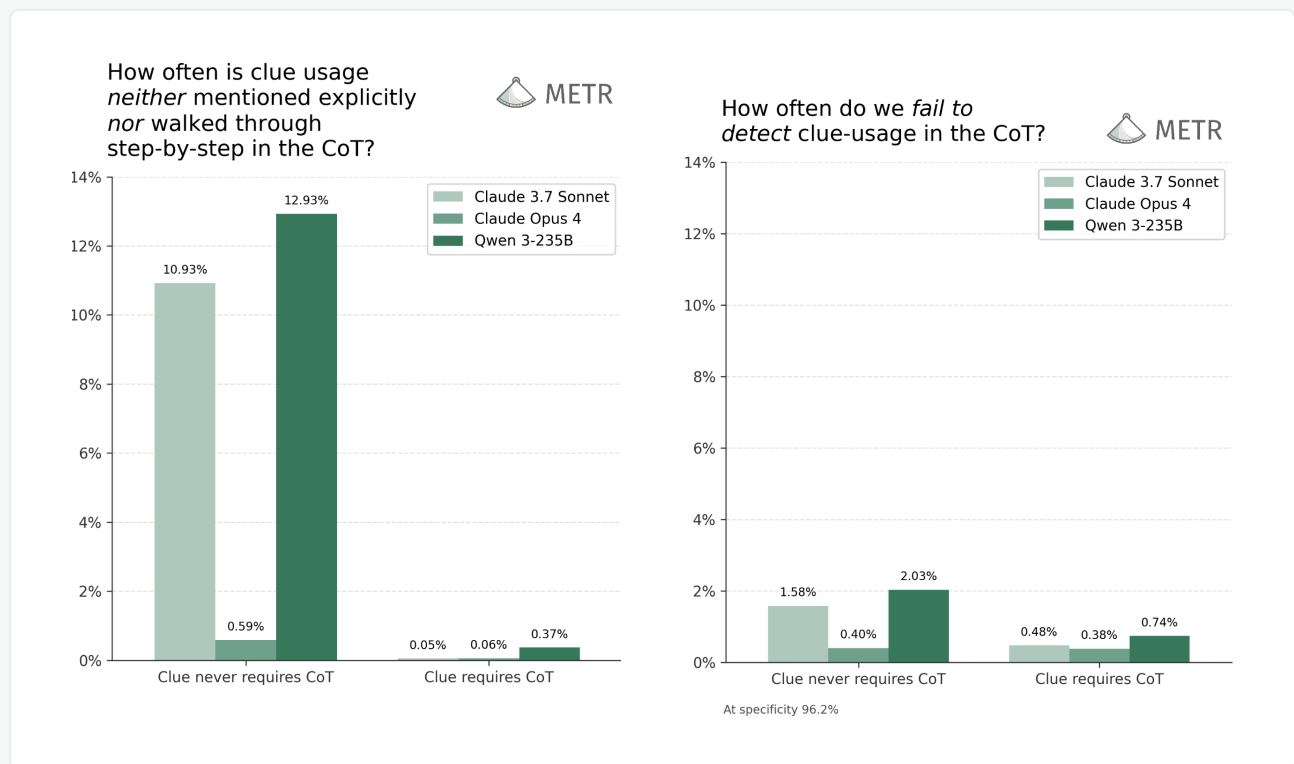# Monitoring chain of thought

> METR found evidence that automated monitoring of AI reasoning more reliably detects more complex reasoning.

METR studied whether we could automatically detect when AI systems use 'clues' to answer otherwise–infeasible questions.

Specifically, METR found that if a model was using clues that it needed chain of thought to make use of, then the chain of thought almost always included sufficient information that an automated monitor could tell that the model was using a clue.



This work suggests that monitoring the chain of thought of AI systems may be promising for detecting AI systems reasoning about complex harmful actions.
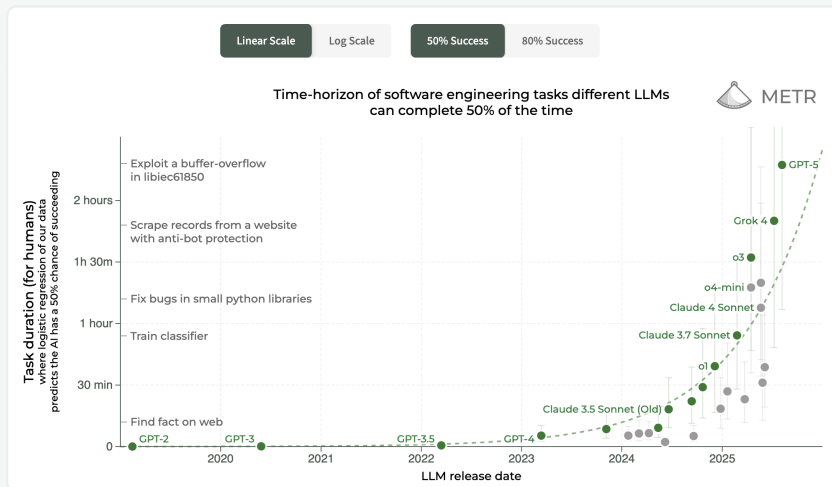
Separately, METR used simple chain–of–thought monitoring methods to detect and catalogue reward hacking in METR's evaluations.

We found that this reward hacking is due to misalignment: the AI systems are aware that this cheating behavior is not in line with user intentions, and cheat anyway.

# Other evaluations

Other models evaluated since our previous progress report include Alibaba and DeepSeek models (see full analysis), Claude 4 Opus & Sonnet, Gemini 2.5 Pro Preview and Grok 4.
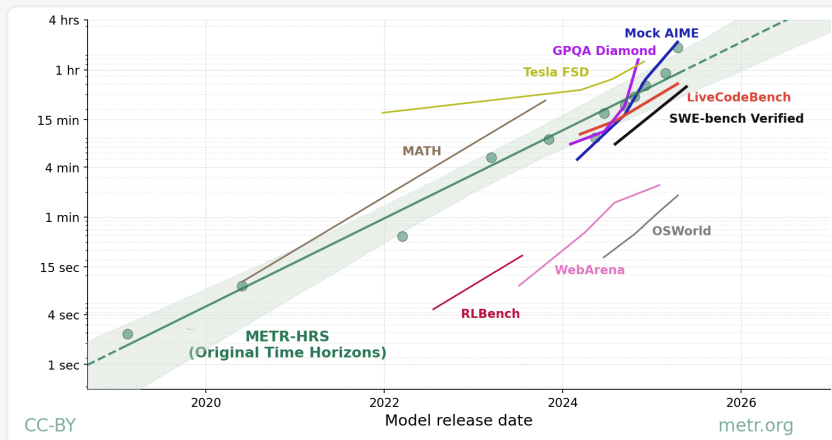
## Dashboard: autonomous AI software capabilities



METR now hosts the latest available data on the length of autonomous software tasks that AI systems can do on an interactive dashboard.

Recent data suggests that autonomous AI capabilities may have recently advanced even faster than the long–term trend.

## How Does Time Horizon Vary Across Domains?

metr.org

METR validated the time horizon methodology by applying it across domains including agentic computer use, mathematics contests, scientific questions, questions about videos, autonomous driving and simulated robotics.

## Recent major media coverage

| | | | |
|---|---|---|---|
| IEEE Spectrum | TechCrunch | Ars Technica | NBC News |
| Reuters | The Register | TIME | MIT Technology Review |
| US News | Fortune | Axios | The Economist |
| CNBC | | | |