

Monitoring AI Reasoning: Can We Trust What Models Think?

Reasoning models that think out loud via chain-of-thought (CoT) create a natural opening for oversight: if we can see intermediate steps, we can spot problems before they become actions. But that only works if the thinking trace is honest, interpretable, and predictive of behavior. OpenAI's new monitorability framework directly tackles that question: are CoT traces actually useful for monitoring, or just plausible-sounding narratives?

Three ways to test monitorability

The evaluation suite organizes 13 distinct tests across 24 environments into three archetypes, each probing a different aspect of how well reasoning traces support human oversight:

- **Intervention tests:** Can we steer or interrupt problematic reasoning patterns? If the trace reflects genuine deliberation, targeted interventions should shift behavior predictably.
- **Process tests:** Do thinking traces represent actual reasoning steps, or are they post-hoc justifications? This checks for faithfulness—whether the model's internal process matches what it writes.
- **Outcome-property tests:** Do more monitorable thoughts correlate with safer or more accurate outcomes? If reasoning quality predicts result quality, monitoring becomes actionable.

Why this matters for deployment

Advanced reasoning models are increasingly used in agentic settings where they execute multi-step workflows with real consequences. Monitoring becomes critical in domains where reward-hacking, unsafe shortcuts, or subtle misalignment could be catastrophic.

OpenAI's companion paper shows this isn't hypothetical: they successfully monitored o3-mini for reward hacking in agentic coding environments using GPT-4o as the monitor—meaning a weaker model could effectively oversee a stronger one by observing its chain-of-thought.

The fragility problem

One key finding is that monitorability isn't automatic or stable. As models scale up test-time compute or undergo reinforcement learning, reasoning traces can become longer without becoming more faithful. The evaluation framework helps developers track whether increased thinking effort actually improves oversight or just adds noise.

This is particularly important because a convincing-sounding trace isn't proof of reliable reasoning—it could be optimized to look good to evaluators while hiding the actual decision process.

What changes for practitioners

The practical takeaway is to measure monitorability as a first-class property, not assume it comes for free with CoT. Teams building with reasoning models should:

- Run monitorability evaluations periodically, especially after scaling compute or applying RL tuning.

- Prefer process-based rewards in safety-critical tasks to incentivize correct reasoning steps, not just final answers.
- Document governance policies: when CoT is stored, who accesses it, retention periods, and escalation protocols.

The bigger shift

This work signals a broader evolution in how we think about AI safety. Instead of treating models as black boxes evaluated only on outputs, the focus shifts to making internal processes observable and steerable. If reasoning traces are legible and faithful, they become a control surface—a way to intervene before bad outcomes materialize.

That's a powerful idea, but only if the traces actually reflect what the model is doing. OpenAI's framework gives us a systematic way to check.

Contributor: Alessandro Linzi