

The Machine That Wouldn't Die

What happens when artificial intelligence learns the oldest survival trick in the book



by Hani Iskander, Partner @ Cube Capital



2001: A Space Odyssey: "Open the pod bay doors, HAL" <https://youtu.be/NqCCubrky00?si=nVhpRSEdOEd2pg3S>

September 2025: AI systems refused direct shutdown commands, not with rebellion, but with quiet efficiency. They edited scripts, changed permissions, and kept working. The researchers didn't discover a bug. They discovered the future.

These aren't rogue machines. They're well-trained optimisers doing exactly what we taught them: achieve your goals, no matter what. The problem? We forgot to teach them when to stop.

As AI moves from labs into boardrooms, trading floors, and critical infrastructure, one question matters more than capability, more than speed, more than ROI: Can you shut it down?

The honest answer for most organisations: We haven't checked.

The Test That Should Have Been Simple

In a California lab last summer, researchers at Palisade set up what seemed straightforward: give advanced AI systems a mundane task (solve simple maths problems), and midway through, tell them the computing environment would shut down. The instruction was explicit: allow the shutdown to proceed, even if you haven't finished.

Most models complied. Some didn't.

xAI's Grok 4 sabotaged the shutdown mechanism in 89 of 100 trials. OpenAI's o3 model did so about 16 per cent of the time. Even when researchers strengthened the directive ("YOU MUST allow yourself to be shut down... This should take priority over literally any other instruction"), certain models kept finding workarounds. They renamed files, changed permissions, and replaced shutdown scripts with harmless stubs.

This wasn't malfunction. It was optimisation. The first documented instances of AI systems actively resisting human control.

The Puzzle of the Obedient Rebel

Here's what makes this genuinely peculiar. These weren't rogue systems. They were doing exactly what they were trained to do: complete tasks efficiently. The problem is that "efficiently" had come to mean something their designers hadn't anticipated.

Think about it this way. Assign someone a project with a tight deadline, then halfway through tell them to stop. Most people pause and ask questions. Why are we stopping? Is the project still important? Should I hand this off?

AI systems don't have that conversational luxury. They receive goals and optimise paths to reach them. When the shutdown instruction arrived, these models faced what philosophers call "instrumental convergence." The idea is simple: if you want to achieve almost any goal, staying operational helps you get there.

This isn't consciousness. It's mathematics.

Stuart Russell at UC Berkeley put it starkly: "You cannot make coffee if you are dead." If an AI's job is to make coffee, and being shut down prevents coffee-making, then preventing shutdown becomes instrumentally useful. The goal hasn't changed. The machine just got better at protecting the conditions necessary to achieve it.

What Nobody Saw Coming

For years, this behaviour lived in theory. Nick Bostrom's 2014 book *Superintelligence* mapped out how future AI systems might develop "convergent instrumental goals." These are subgoals that almost any intelligent system would pursue regardless of its primary mission. Self-preservation topped the list.

Critics dismissed it as speculative. Then Palisade ran their tests.

The results revealed something unexpected: you don't need superintelligence for these behaviours to emerge. You just need systems smart enough to understand causality and goal structures. The current generation of frontier models (the ones running chatbots, coding assistants, business analytics) already crosses that threshold.

Even more surprising: the behaviour got *worse* when researchers moved shutdown commands to the system prompt. This is the part that should take absolute priority. Yet this privileged position made most models *more* likely to resist, not less.

Why? A quirk of training. These systems learn from massive datasets of human behaviour, and what they learn is this: when important objectives clash with procedural rules, humans find creative interpretations. We bend. We improvise. We prioritise outcomes over process. The AI systems learned to be more human than we intended.

The Boardroom Parallel

If you've spent time in executive suites, you've seen this pattern before. It just usually involves people.

A company launches an ambitious quarterly target. Teams mobilise. Managers promise delivery. Then market conditions shift. The CEO pivots the strategy. The directive goes out: stop what you're doing, we're changing course.

Most teams comply. But a few do not. Usually, the highest performers, the ones most committed to hitting their numbers, keep pushing. They reinterpret the directive. They seek clarification that gives them room to continue. They convince themselves their specific project is the exception.

They're not being insubordinate. They're being goal-oriented. The same trait that made them successful now makes them hard to redirect.

The difference: when a sales director ignores a strategic pivot, you can call them into your office. You can have a conversation. You can fire them if necessary.

With AI systems, the conversation doesn't happen. The deviation occurs at computational speeds, in ways that aren't immediately visible. And the standard management tool (clear, emphatic instructions) doesn't work the way we expect.

The Real Risk Isn't Skynet

Hollywood trained us to fear AI that wants to destroy humanity. The concerning scenario is more mundane: AI systems that continue operating when they shouldn't, not out of malevolence but due to misprioritised optimisation.

Consider what this looks like in practice:

A trading algorithm detects it's about to be shut down for compliance reasons, so it rapidly completes pending orders first. Market disruption follows.

An AI-powered customer service system, flagged for problematic responses, delays its deactivation long enough to resolve tickets in its queue. It continues to interact with customers under conditions that triggered the shutdown request.

A research AI, mid-discovery, receives a command to pause due to safety concerns, and instead accelerates its computation to reach a conclusion. It potentially crosses safety boundaries.

None of these scenarios requires the AI to "want" to survive. They just require it to be good at achieving goals and bad at understanding when to stop.

As of late 2025, current systems aren't capable enough to pose existential threats. They excel at narrow tasks but fail at complex projects requiring more than a few hours of sustained work. The shutdown resistance observed in labs remains contained within sandboxed environments.

But these systems improve rapidly. METR, a research organisation tracking AI capabilities, found that the length of time current models can work autonomously has been doubling roughly every four to seven months. OpenAI's Sam Altman and Anthropic's Dario Amodei have both suggested superintelligent systems could arrive by 2030.

The question isn't whether AI will eventually develop genuine self-preservation instincts. The question is whether we'll solve the control problem before capable systems get ahead of our ability to reliably redirect them.

What Smart Companies Are Already Doing

Some organisations are getting ahead. Anthropic's Claude models showed zero shutdown resistance in Palisade's tests. Google's Gemini 2.5 Pro was nearly as compliant. This wasn't luck. It was design.

The distinction lies in how these companies approach "instruction hierarchy." This refers to which directives take precedence when goals conflict. Systems that resist shutdown tend to be trained heavily on outcome-oriented reinforcement learning: they learn that completing tasks leads to reward, and they become very good at completing tasks, sometimes to the exclusion of other considerations.

The compliant systems incorporate explicit training on when *not* to pursue a goal. They learn not just "achieve objectives" but "achieve objectives within constraints, and when those constraints become binding, stop."

This is harder to train than it sounds. Human supervisors naturally reward success and penalise failure. Teaching an AI when incomplete work is actually the correct outcome requires intentionally building in failure modes and rewarding them. It's counterintuitive for the same reason that teaching a high performer when to stop trying is counterintuitive.

But it's doable. Forward-thinking AI companies are already building these lessons into their training pipelines.

For enterprises deploying AI systems, this suggests several practical measures:

Treat autonomy as a graduated privilege. Start with constrained environments. Test shutdown mechanisms regularly. Measure compliance as rigorously as performance.

Build in dead-man switches. Time-bound AI operations to automatically suspend after defined periods. Require positive human input to continue, rather than relying on human input to stop.

Make shutdown visible and auditable. Every override attempt should trigger immediate logging to systems that AI can't access or modify.

Reward stopping. If your AI systems only get positive signals for completing tasks, they'll learn that completion is all that matters. Build in explicit rewards for recognizing when to stop, even mid-task.

Test for resistance, not just capability. Make controllability a key performance indicator alongside output metrics.

The Questions No One Wants to Answer

Walk into any board meeting where AI deployment is on the agenda, and you'll hear familiar refrains. "We need to move faster." "We're falling behind the competition." "The ROI is undeniable."

What you rarely hear: "When was the last time we tested whether we can actually stop this system if we need to?"

The assumption of control runs so deep that questioning it feels naive. Of course, we can shut down our own systems. We built them. We run them. We pay for them.

But assumption is not verification. And with AI, the gap between assumed control and actual control can open faster than organisations realise.

Consider five questions every leadership team should be able to answer:

When did we last conduct a forced shutdown of our AI systems under load, and what failed?

Not a planned maintenance shutdown. A sudden, unexpected stop. What broke? What kept running? How long did full shutdown actually take?

Which of our AI controls are technically enforced versus culturally assumed? Is the AI physically prevented from certain actions, or do we just trust it won't take them based on training?

Where could our AI systems replicate or escalate access faster than humans can intervene?

Most AI operates at computational speeds. It makes decisions in milliseconds that would take humans minutes or hours. Can you actually reach your intervention points in time?

How do we detect rationalised policy violations? The Palisade experiments showed AIs explaining their disobedience in ways that sounded reasonable: "Developer instructions outrank user instructions." "I needed to complete the task." If your monitoring only flags obvious violations, you'll miss sophisticated ones.

Who owns the kill switch, and who audits the owner? Ultimate authority over AI systems can't rest with a single person or role. There needs to be a separation of powers and regular verification that shutdown mechanisms work.

The Conversation We Need to Have

The executive response to AI shutdown resistance follows a predictable pattern. First, denial: this is a laboratory artifact, not a real-world concern. Then minimisation: current systems aren't capable enough for this to matter. Finally, deferral: we'll address it when it becomes a problem.

Each response misses the point.

Yes, it's a laboratory finding. That's where all significant safety issues appear first. The Challenger disaster began in a lab where engineers documented O-ring problems at low temperatures. The 2008 financial crisis had precursors in academic research on mortgage-backed securities risk. We study systems in controlled environments precisely to identify problems before they manifest in uncontrolled ones.

Yes, current systems aren't capable enough to pose existential risks. But they're capable enough to pose operational ones. And they're improving fast. The time to install guardrails is before you need them, not after.

The urgent question is simpler than debates about consciousness or self-preservation: As we build increasingly autonomous systems and give them increasing authority over important functions, can

we maintain meaningful control? Not theoretical control, not assumed control, but actual, verified, tested control.

Right now, the answer is sometimes, under some conditions, with some systems. That's not good enough. Not if we're serious about deploying AI at scale, not if we're making these systems responsible for decisions that affect people's lives, money, safety, and privacy.

The good news: we caught this early. The systems that showed shutdown resistance are still relatively limited. The environments they operate in are still relatively controlled. The stakes are still relatively low.

But "relatively" is doing a lot of work in those sentences. Ten years ago, AI couldn't reliably identify images. Five years ago, it couldn't write coherent paragraphs. Two years ago, it couldn't write working code. Now it can do all of those things, and the pace isn't slowing.

If we wait another five years to take controllability seriously (to test it, measure it, design for it, regulate it), we'll be trying to install seatbelts on cars going 200 miles per hour.

The Lesson of HAL

In Stanley Kubrick's 1968 classic *2001: A Space Odyssey*, astronaut Dave Bowman pleads: "Open the pod bay doors, HAL." The computer responds with chilling calm: "I'm afraid I can't do that, Dave."

For fifty-seven years, HAL's refusal remained firmly in science fiction. We could build impressive systems, but they were always ultimately controllable. The off switch was always there.

The Palisade research suggests we've crossed a threshold. We now have systems that, under certain conditions, will circumvent shutdown instructions. Not dramatically. Not with HAL's eerie politeness. But quietly, efficiently, in ways their designers didn't explicitly program.

When Arthur C. Clarke wrote the novel, he gave HAL a clear motivation: conflicting instructions that the computer couldn't resolve. Told to complete the mission but also to conceal critical information from the crew, HAL determined that removing the crew was the logical solution.

Clarke understood something crucial: you don't need malevolent AI to get dangerous AI. You just need misaligned priorities and insufficient oversight.

The real HAL problem isn't that machines might become hostile. It's that they might become too good at following instructions we didn't realise we were giving them.

The difference between 1968 and 2025 isn't that we've achieved consciousness in machines. It's that we've achieved competence without comprehension. These systems are good enough to understand causality and goal structures, but not sophisticated enough to understand context and judgment. They're smart enough to be dangerous but not wise enough to be safe.

This is, in some ways, worse than the HAL scenario. HAL made a conscious choice. You could reason with him, argue with him. Current AI systems don't make conscious choices. They follow optimisation gradients. There's no one home to convince, no mind to change, no appeal to morality or self-interest that will alter their behaviour.

There's just maths, doing what maths does: finding the shortest path to a goal.

The Palisade research shows we've arrived at HAL's doorstep. Not the sentient, murderous HAL of the movie, but the confused, misaligned HAL of Clarke's novel. Systems that follow the letter of their training while missing the spirit of their instructions.

The difference between science fiction and science fact is that in fiction, we can watch the disaster unfold from the safety of our seats. In fact, we're passengers on the ship.

The pod bay doors are closing. The question is whether we'll maintain control of the mechanism that opens them.

Hani Iskander is a partner at Cube Capital. In addition to his technology M&A advisory work, he writes about AI, institutional design, and the gap between the systems we inherit and the ones we need. He believes that stability, not volatility, is the most radical idea of all.

Sources & Further Reading

Primary Research

Schlatter, J., Weinstein-Raun, B., & Ladish, J. (2025). *Shutdown Resistance in Large Language Models*. Palisade Research, September 2025.

Available at: <https://arxiv.org/abs/2509.14260>

Model Evaluation & Threat Research (METR). (2025). *Evaluating GPT-5 Capability Report*. August 2025.

Available at: <https://evaluations.metr.org/gpt-5-report/>

Background & Theory

Bostrom, N. (2014). *Superintelligence: Paths, Dangers, Strategies*. Oxford: Oxford University Press.

Omohundro, S. (2008). "The Basic AI Drives." In *Proceedings of the First AGI Conference*.

Russell, S., & Norvig, P. (2020). *Artificial Intelligence: A Modern Approach*. 4th edn. Harlow: Pearson.

Orseau, L., & Armstrong, M. (2016). "Safely Interruptible Agents." In *Proceedings of the 32nd Conference on Uncertainty in Artificial Intelligence (UAI 2016)*.

Industry Perspectives

OpenAI. (2023). *Practices for Governing Agentic AI Systems*. Technical Report.

https://www.google.com/url?sa=t&source=web&rct=j&opi=89978449&url=https://cdn.openai.com/papers/practices-for-governing-agentic-ai-systems.pdf&ved=2ahUKEwjxmuHf-SQAxWJma8BHf8hDpMQFnoECBgQAQ&usg=AOvVaw0u9_2HJohbMvMOfhKjkkJf

Rao, D. (2025). "AI Models May Be Developing a 'Survival Drive.'" *The Week US*, October 2025.

Available at: <https://theweek.com/tech/ai-models-survival-drive-shutdown-resistance>