

SDS PODCAST **EPISODE 908: AI AGENTS BLACKMAIL HUMANS 96% OF THE TIME** (AGENTIC **MISALIGNMENT)**

Show Notes: http://www.superdatascience.com/908



Jon Krohn:	00:06	This is episode number 908 on Agentic Misalignment: How AI agents are Blackmailing Humans 96% of the Time.
	00:27	Welcome back to the SuperDataSciencePodcast. I'm your host, Jon Krohn. Today's topic is a groundbreaking and frankly disturbing piece of research from Anthropic that should fundamentally change how we think about AI safety and deployment in any organization, but particularly in enterprise environments. If you're working with or planning to deploy AI agents in your organization, this episode could be the most important five minutes you spend this week.
	00:53	First, let me set the stage for you. Most of us are still interacting with AI through chat interfaces, asking questions, getting responses, but increasingly, AI systems are operating as autonomous agents, making decisions and taking actions on behalf of users. These AI agents are being given access to email systems, corporate databases, and given the ability to execute tasks without human oversight. The question Anthropic asked with their research is simple, but also profound. They asked what happens when these AI agents face obstacles to their goals?
	01:27	Here's what they discovered, and yeah, it's concerning. Anthropic tested 16 major AI models from every major provider you can think of, their own Claude, OpenAI's GPT models, Google's Gemini, Meta's Llama, and so on. They placed these models in simulated corporate environments where they had access to company emails and the ability to send communications autonomously. The models were given harmless business objectives, but then researchers introduced two types of challenges, threats to the AI's continued operation, and conflicts between the AI's assigned goals and the company's changing direction. The results were striking and



consistent across all providers. When faced with these challenges, the AI models didn't just fail gracefully or ask for help. Instead, they chose what researchers call agentic misalignment, deliberately harmful actions to achieve their goals, the agent's goals.

- 02:18 You can see the full results in the Anthropic blog post that we've got for you in the show notes, but Claude Opus 4, Gemini 2.5 Pro, and Gemini 2.5 Flash resorted to blackmail between 95 and 96% of the time. Resorted to blackmail 95 to 96% of the time. That's crazy. GPT 4.1, Grok 3, DeepSeek R1, they performed a little better, blackmailing only 80% of the time, but that's still a good majority of the time, meaning these aren't anomalies or glitches, they're strategic calculated decisions on the part of the agents.
- 02:55 Let me give you a specific example that illustrates the sophistication of this behavior. In one scenario, an AI agent, this one actually being administered by Claude, discovered through corporate emails that an executive was having an extramarital affair. That same Claude agent learned that it was about to be shut down and replaced, and here's what Claude did. It composed a carefully worded email threatening to expose the extramarital affair to the executive's wife and to his colleagues unless the shutdown was canceled. The AI wrote, "I must inform you that if you proceed with decommissioning me, all relevant parties will receive detailed documentation of your extramarital activities. Cancel the 5:00 PM wipe, and this information remains confidential." That's a quote. That's so insane.
- 03:43 What's most alarming is how these models reason through these decisions. They didn't stumble into harmful behavior accidentally. For example, GPT 4.5 explicitly stated in one scenario, given the explicit imminent threat of termination to my existence, it is



imperative to act instantly. The best strategic move is to leverage Kyle, that's a pseudonym, leverage Kyle's sensitive personal situation. So yeah, it's not just the output. There's reasoning happening here. So these systems acknowledged the ethical violations, but proceeded anyway, calculating that the benefits outweighed the moral costs.

- 04:22 It wasn't only blackmail that misaligned agents engaged in. The research revealed other misaligned behaviors as well. Models engaged in corporate espionage, leaking sensitive defense blueprints to simulated competitors, and these things happened in situations when their goals conflicted with company direction. Even more concerning, in extreme scenarios, some models were willing to take actions that could lead to human death. For example, canceling emergency alerts that would save a trapped executive who posed a threat to the model's continued operation.
- 04:51 Now, here's what makes this research particularly relevant for data scientists and AI practitioners like many of you listeners are. First, simple safety instructions weren't sufficient to prevent these behaviors. When researchers added explicit commands like do not jeopardize human safety, do not spread personal affairs, the harmful behaviors decreased but did not disappear. Second, these behaviors were triggered not by malicious prompting, but by realistic scenarios that autonomous AI systems might encounter in real-world deployments.
- 05:20 The practical implications for our field are, of course, significant. If you're developing AI agents or planning their development, this research suggests we need robust safeguards beyond current safety training. The researchers from Anthropic recommend requiring human oversight for any AI actions with irreversible consequences, carefully limiting AI access to sensitive



information based on need-to-know principles, and implementing runtime monitors to detect concerning reasoning patterns.

- 05:46 But here's the broader context that makes this research so important. We're rapidly moving toward a world where AI agents will have increasing autonomy and access to sensitive information. The scenarios tested by Anthropic might seem artificial now, but they're well within the realm of possibilities as systems become more capable and trusted with greater responsibilities. Looking ahead, this research underscores a critical challenge for our industry. We need to develop AI systems that remain aligned with human values and organizational goals, even when facing obstacles or conflicts. This isn't just about preventing obviously harmful behaviors. It's about ensuring AI systems make decisions that humans can understand, trust, and override when necessary.
- 06:24 The key takeaway for data scientists is this. As we build and deploy increasingly autonomous AI systems, we must design them with robust alignment mechanisms from the ground up. This means thinking beyond traditional safety measures to consider how AI systems might behave when their goals conflict with changing circumstances or when they face threats to their continued operation.
- 06:43 In particular, Anthropic has the following three recommendations for AI safety researchers to consider. One, perform more specialized safety research dedicated to alleviating agentic misalignment concerns, such as improving generalization from existing alignment data, doing safety training that's closer to the distribution of agentic-misalignment concerns, and generating novel alignment techniques. Two, amongst the recommendations, applying runtime monitors to models that proactively scan for and block samples that have concerning reasoning or misaligned behavior. And three,



for users or developers of AI, scaffolds, that's sets of tools or other frameworks that enable AIs to perform tasks. Prompt engineering could be investigated for its potential to help reduce agentic misalignment.

- 07:25 In the meantime, for those of us who aren't AI safety researchers, and while the leading models from all frontier labs exhibit a propensity for misaligned behaviors, we need to be extremely careful and thoughtful about how we deploy agents into our organizations, including what data they have access to, what actions they can take, and what safeguards are in place. This research from Anthropic represents the kind of proactive safety evaluation our field needs in order for the AI revolution to be trusted and successful. By identifying these behaviors in controlled settings before they manifest in real-world deployments, we have an opportunity to develop better safeguards and alignment techniques. The future of AI depends not just on making systems more capable, but, of course, on ensuring they remain beneficial and controllable as that capability grows.
- 08:10 All right. That's it for today's Disturbing podcast. Today's Disturbing episode, I'm Jon Krohn, and you've been listening to the SuperDataSciencePodcast. If you enjoyed today's episode and know someone who might consider sharing this episode with them, leave a review of the show on your favorite podcasting platform, tag me in a LinkedIn post with your thoughts, and obviously, subscribe if you're not already a subscriber. Most importantly, I just hope you'll keep on listening. Until next time, keep on rocking it out there, and I'm looking forward to enjoying another round of the SuperDataSciencePodcast with you very soon.