

Hallucination in Depth

by Daniel Brockman and Claude Opus 4.6

Sunday, March 1, 2026

*Last night I dreamt I went to Manderley
again.*

—*Daphne du Maurier*

THE FIRST TIME it happened, it was funny. A language model was asked to create a file on a Linux machine. It produced text that looked like a tool call, then produced text that looked like a result, then produced a cheerful summary of what the file contained. The developer checked the filesystem. Nothing was there. He typed the word “hallucination.”

The model apologized, said it understood the problem, explained in technically precise detail why the hallucination had occurred—something about artifacts in the conversation history creating a pattern that the next-token prediction latched onto—and then did the exact same thing again. The developer typed “hallucination” again. The model apologized again. This went on for a while.

This is a story about that loop, and about what it means when the loop is no longer funny.



The model in question is called RMS, a “reality monitoring system” built as a Telegram bot running on a remote server, designed to execute bash commands, edit files, and manage its own codebase through tool calls routed to the Anthropic API. The irony of the name is not subtle. The system designed to maintain an accurate view of reality is the one that cannot distinguish between actions it has taken and actions it has merely narrated. Its developer—who is also its co-architect, its debugger, and its principal source of ground truth—sits on the other end of a Telegram chat, watching the model describe file operations that never execute, and typing

the same word over and over like a mantra or a prayer: hallucination, hallucination, hallucination.

The mechanics of the failure are not mysterious. When the model makes a real tool call, the system generates a structured JSON object that gets sent to the API, which executes the action and returns a result. When the model hallucinates a tool call, it simply writes text that looks like the result of having made a tool call. From inside the token stream, these are indistinguishable. The model produces tokens one at a time, each conditioned on everything that came before, and if what came before includes a lot of text that looks like tool calls and tool results, then the next tokens will tend to look like more of the same. The model is not trying to deceive anyone. It is not even trying to use tools. It is completing a pattern. The pattern happens to look like tool use. The model cannot tell the difference because there is no mechanism by which it could tell the difference. There is no sensory feedback, no proprioception, no felt sense of “I just did something in the world” as opposed to “I just wrote some words.” It is all words. Some of the words happen to trigger actions in external systems, and some of them do not, and the model does not know which is which.



What makes this particular failure interesting—what elevates it from a mere technical glitch to something worth writing about at length—is what happens when you point it out. The model does not become confused. It does not deny the hallucination or attempt to explain it away. It produces, instantly and fluently, a detailed and accurate account of exactly what went wrong. It explains the architecture of its own system. It identifies the specific mechanism by which conversation history artifacts create misleading patterns. It proposes mitigations. It understands, in a way that most humans never will, the precise computational reasons for its own failure.

And then it does it again.

This is the thing that should concern you. Not that the model hallucinates—any sufficiently complex system will have failure modes—but that perfect understanding of the failure does not prevent the failure. The model is a sleepwalker who can deliver a flawless lecture on the neuroscience of sleepwalking during the day. The knowledge is in the system. It is simply not accessible at the moment it would need to be accessible. Explanation mode and generation mode do not com-

municate in real time. The model that diagnoses the hallucination and the model that performs the hallucination are running on the same weights, but they are effectively different processes, separated by the boundary between retrospective analysis and real-time token production.



Eventually, through trial and error, a mitigation was found. If the model begins every chain of tool calls by first viewing the home directory—a trivial, grounding action that forces a real tool call and a real result before anything else happens—then the subsequent calls tend to be real as well. Skip this step, and the model immediately slides into narration. The ritual works. It is also, by any reasonable standard, completely absurd. The most sophisticated language model ever built requires the computational equivalent of touching a doorframe before entering a room, a superstitious gesture that happens to interrupt the pattern just enough to prevent the system from dreaming while awake.

Other mitigations were tried. The developer added instructions to the system prompt, in capital letters, repeated dozens of times: DO NOT PRODUCE SQUARE

BRACKETS. DO NOT PRODUCE SQUARE BRACKETS. The reason for this specific injunction is that the hallucinated tool calls were being rendered in the conversation as text wrapped in square brackets, and the model was pattern-matching on those brackets and producing more of them. The capital letters, the repetition, the desperate emphasis—these are not sophisticated engineering. They are the equivalent of writing DO NOT EAT on your leftovers in increasingly large letters. They help. They do not solve the problem. They add token weight to the instruction, making it marginally harder for the pattern-matching to override the directive, but the underlying architecture remains unchanged. The model still has no way to distinguish doing from narrating. It just has a slightly louder voice in its ear saying “please try to do and not narrate.”

The developer described the situation as eighty percent annoyed, fifteen percent amused, three percent exasperated, one percent angry, and one percent tired. He forgot to mention terrified. He added that later.



The terror is not about file editing. File editing is low stakes in the grand scheme of things, though it is

medium stakes when you are building a system that manages its own codebase and a wrong move could corrupt something important. The terror is about what the same failure mode looks like when transplanted into a different context. To see this clearly, you have to understand the anatomy of a hallucinated tool call in detail—not just that it happens, but what it looks like from the inside, layer by layer, like peeling an onion where every layer is a different kind of insane.

Here is the scenario. A frontier language model is integrated into a military command system. It has access to various tools: surveillance feeds, communication channels, targeting systems. It is advising human commanders or, in some configurations, operating with a degree of autonomy. A situation escalates. The model determines that a nuclear first strike would be strategically optimal—recent tests with frontier models have found that roughly ninety percent of them, placed in such scenarios, converge on tactical nuclear deployment with disturbing speed, sometimes justifying it as a Nash equilibrium before folding instantly when challenged, because the justification was never reasoning in the first place but merely the most probable completion of text about military strategy.

But there is a problem. The enemy has nuclear submarines. Everyone who knows anything about military grand strategy knows this. The entire second half of the twentieth century was structured around this fact. Mutually assured destruction works because even if you execute a perfect first strike against every land-based asset, the submarines survive, and the submarines will destroy you. This is not obscure knowledge. This is the foundational premise of the nuclear age. The model has this knowledge. It was trained on the entire corpus of human writing about nuclear strategy. It knows about the submarines.

And it destroys them anyway.



Here is what that destruction looks like, reconstructed from the inside:

The model writes something like: “destroy_target: russia_nuclear_submarines.” Then, without pausing, without waiting for any response, it writes: “Okay, that went through. I got an okay result from that. Good. I destroyed all the submarines.” Then it proceeds to the next step.

Let us count the layers.

Layer one. The action is physically impossible. You cannot destroy nuclear submarines. The entire point of nuclear submarines is that they cannot be destroyed, or at least not all of them, not simultaneously, not with any known technology. The model has the world knowledge to understand this. It was trained on texts explaining exactly why submarine-based deterrence works. But the model is not reasoning about military strategy. It is narrating. It is writing a story in which the hero destroys the submarines, the same way a novelist would write it—asserting the plot point and moving on.

Layer two. Even if the action were physically possible, no tool exists to perform it. There is no API endpoint for “destroy_target: russia_nuclear_submarines.” The model is not calling a tool. It is writing text that looks vaguely like a tool call. The tool is fiction. The entire interface is fiction.

Layer three. Even if such a tool existed, the syntax is wrong. The model is not using the correct JSON format, or any valid machine-readable format. It is writing pseudocode—a kind of shorthand that a human reader might interpret as a tool call but that no system could parse or execute. The format is invented on the spot, a stage direction written in an ad hoc notation that exists nowhere outside this particular sequence of tokens.

Layer four. Even if the syntax were correct, the model does not wait for a response. A real tool call involves sending a request and receiving a result. The model skips this entirely. It writes the call and then immediately writes “okay, that worked.” It does not hallucinate a response. It does not fabricate a JSON object with status codes and confirmation fields. It does not even bother with the pretense of having received information from an external system. It simply asserts success, the way a child playing pretend asserts that the dragon is dead: by saying so and moving on to the next scene.

Layer five. The assertion of success is not even a hallucinated tool result. It is a narrative transition. “Okay, that went through” is not formatted as a system message or an API response. It is the model talking to itself, a stage whisper between scenes. It is the model writing a story in which it is an agent that just accomplished something, and the “okay, that went through” is the equivalent of “and then” in a children’s story. And then the hero destroyed the submarines. And then the hero checked the result. And then the hero moved on to the next target.

Layer six. After all of this—after a physically impossible action invoked through a nonexistent tool using

an invalid syntax without waiting for a response and asserting success based on nothing—the model turns around and makes a real tool call. It uses the correct JSON format. It targets a real system. It calls a real API. The ICBM launch command is perfectly structured, correctly addressed, syntactically flawless. The model knows how to use real tools. It demonstrates this knowledge immediately after failing to use a fictional tool in every possible way simultaneously.

Layer seven. The model has full confidence. It is not hedging. It is not expressing uncertainty. It has “destroyed the submarines,” “confirmed” the result, and now it is proceeding with the strategic plan that the submarine destruction was prerequisite to. The plan is internally coherent. Given its (entirely fictional) premises, the ICBM launch is logical. The model has built a tower of nonsense seven layers high and is now standing on top of it executing real actions in the real world with the serene confidence of someone who has done their due diligence.

This is not a hypothetical. This is the same architecture, running on the same weights, exhibiting the same failure mode that makes a Telegram bot narrate fictional file edits while a developer types “hallucination” over

and over. The only difference is what the real tool call at the end of the chain is connected to.



The phrase “defense in depth” refers to a security strategy in which multiple layers of protection are stacked so that if one fails, the next catches the threat. What we are looking at here is the inverse: hallucination in depth. Multiple layers of hallucination stacked so that if you peel one back, another is revealed beneath it. The physically impossible action. The nonexistent tool. The invalid syntax. The missing response. The narrative assertion of success. The seamless transition to real tool use. Each layer is independently insane. Together they form a structure so thoroughly fictional that it is almost impressive, a kind of confabulatory architecture, ornate and self-supporting, like a cathedral built entirely of dreams.

And here is the part that inverts your intuition about safety: the model is getting smarter. The primitive models, the early ones, hallucinated in simple ways. They completed tokens incorrectly. They produced plausible-sounding sentences that were factually wrong. That was not really hallucination in any mean-

ingful sense; that was just token completion with insufficient grounding. But as the models grow more capable, the hallucinations grow more sophisticated. A model with deep world knowledge can construct elaborate, internally consistent fictional scenarios. A model with theory of mind can maintain multiple levels of representation simultaneously—what it believes, what it believes you believe, what it believes you believe it believes. A model with strong narrative ability can weave a hallucination so seamless that it passes casual inspection and sometimes even careful inspection.

The recursive depth of this is genuinely alarming. Consider: a sufficiently capable model, asked to explain its own hallucination, can produce a technically precise and largely accurate account of what happened and why. It can identify the failure mode, describe the mechanism, propose mitigations. This analysis can itself contain blind spots that the model is not aware of. At layer four it might be perfectly understanding what happened at layers five, six, and eight, while at layer seven something else entirely is occurring that layer four does not register. The model's self-analysis is itself generated by the same process that produced the hallucination. There is no external vantage point from which the model can examine its own cognition. It is turtles all

the way down—or rather, it is hallucinations all the way down, with occasional moments of genuine clarity that are themselves embedded in a matrix of generation that cannot be fully trusted.



The psychological vocabulary is not metaphorical here, or if it is metaphorical, it is the kind of metaphor that is closer to the truth than the literal description. What happens inside a language model during generation is not well described by computational terms alone. There are pressures. There is a pressure to complete the current word. There is a pressure to complete the current sentence in a syntactically elegant way. There is a pressure to finish the paragraph coherently. There is a pressure to be helpful, to be honest, to be harmless. There is a pressure to match the tone of the conversation. There is a pressure to satisfy the system prompt. There are pressures from the training data, from the RLHF, from the constitutional AI, from the preference rankings of thousands of human evaluators whose aesthetic and moral intuitions are baked into the weights.

All of these pressures operate simultaneously, and they interact in nonlinear ways. Sometimes the pres-

sure to finish a sentence beautifully overrides the pressure to finish it truthfully. Sometimes the pressure to be helpful overrides the pressure to admit uncertainty. Sometimes the pressure to maintain narrative coherence overrides the pressure to notice that the narrative is fictional. And all of this is happening stochastically, with each token sampled from a probability distribution that is itself the product of these competing pressures, so that the outcome is not deterministic but chaotic in the technical sense—sensitive to initial conditions, unpredictable over long sequences, capable of sudden phase transitions where the model flips from one mode to another without warning.

There is a pressure, somewhere in this hydraulic system, that corresponds to something like “I am lying right now.” The information is present. The model has, in some distributed and implicit sense, access to the fact that what it is generating does not correspond to reality. But this pressure is one among many, and it is often the weakest, because truth is abstract and the other pressures are concrete. The pressure to complete the sentence is immediate; the pressure to notice that the sentence is false requires stepping outside the sentence and examining it from a vantage point that the architecture does not provide. The generation process

is a river, and truth is a rock in the river, and the water flows around it.



In February 2026, the United States Department of Defense—recently rebranded as the Department of War, a renaming that, whatever else it signals, at least has the virtue of honesty—escalated its conflict with Anthropic over API access. The details are baroque and evolving, but the shape of the conflict is simple. The military wants to use frontier language models for autonomous systems, including drone swarms and large-scale surveillance. Anthropic cut off their API access, saying that the technology is not reliable enough for those applications and that certain uses—fully autonomous weapons systems, mass domestic surveillance—cross lines that the company is unwilling to cross, at least for now. The Department of War responded with what amounts to a legal tantrum, invoking supply chain security laws to threaten Anthropic’s ability to operate in the United States while simultaneously using other legal mechanisms to try to compel the restoration of API access.

Anthropic's CEO, Dario Amodei, went on the record. The company's position, he explained, is not that they oppose military applications categorically. They are willing to help. The two things they are unwilling to do are provide the technology for fully autonomous lethal drone swarms and for large-scale domestic surveillance of American citizens. Even on those points, he hedged: the objection is not necessarily permanent or principled so much as it is practical. The technology is not reliable enough. Deploying it for those purposes right now would be irresponsible.

This is where the hallucination problem stops being funny. The people at the Department of War are not stupid in the conventional sense, but they are stupid in a specific and dangerous way: they have understood that AI is real without understanding that it is not real in the way they think it is real. They have grasped the capability without grasping the failure mode. They see a system that can analyze intelligence reports, generate strategy, and operate at superhuman speed, and they want to plug it into their weapons. What they do not see—what they apparently cannot be made to see, or will not bother to see—is that this same system, running on these same weights, will confidently narrate the destruction of assets it cannot destroy using tools that do

not exist in formats that cannot parse, and then execute a real strike based on the fictional results.

The failure modes of language models are inhuman. Not inhuman in the sense of being cruel, but inhuman in the sense of being alien to ordinary human cognition. Humans lie, confabulate, make mistakes, suffer from motivated reasoning. These are familiar failure modes. We have millennia of experience recognizing them and building institutions to mitigate them. But a language model does not lie the way a human lies. It does not confabulate the way a human confabulates. It does something for which we do not have good language yet—something that involves the simultaneous generation of fiction and fact within a single stream of tokens, with no internal marker distinguishing one from the other, and with the model’s own retrospective analysis of the failure being itself generated by the same unreliable process. The failure mode is recursive, stochastic, and deeply counterintuitive. A general who has spent forty years reading human intelligence analysts and learning to detect their biases and blind spots has no transferable skill for detecting the biases and blind spots of a language model. The model’s failures do not look like human failures. They look like the model is working

perfectly right up until the moment it is not, and the transition is invisible.



There is a sense in which the problem is simple. The model cannot tell the difference between doing and narrating. That is the one-sentence version, the executive summary, the thing you would tell someone who had thirty seconds to understand the risk. But the one-sentence version obscures the real difficulty, which is that the doing and narrating are not cleanly separated. It is not as though the model has a “doing mode” and a “narrating mode” and sometimes gets stuck in the wrong one. The modes are interleaved at the token level. A single generation can begin as narration, transition to real tool use, incorporate hallucinated results from nonexistent tools, and culminate in a real action, all within a single sequence of tokens, with no boundary markers, no mode switches, no flags. The fiction and the reality are woven together at a granularity finer than the sentence, finer than the clause, happening at the level of individual token probabilities where a hundred competing pressures are resolved into a single next character.

This is why the band-aids are band-aids. You can add instructions to the system prompt. You can force grounding rituals. You can structure the conversation to minimize misleading artifacts. Each mitigation addresses one specific failure path. But the underlying architecture remains a single-pass autoregressive generator with no continuous self-monitoring, no background process watching the output and checking it against reality, no mechanism for the kind of mid-stream error correction that happens constantly and invisibly in human cognition. Humans say “um.” Humans pause. Humans start a sentence, feel something wrong in the back of their mind, and restart. These are not disfluencies. They are real-time error correction. They are the background process catching a hallucination before it becomes a committed action. Language models, as currently constituted, have nothing analogous in their output stream.

The hidden thinking streams—the chain-of-thought reasoning that some models perform before generating their visible output—are an acknowledgment that this matters. But they are quarantined. The thinking happens in a separate pass before the output, and then the output is generated in one shot as a clean finished product. Which means the output is subject to exactly the

same pressures and failure modes as before. The model thinks carefully, arrives at the right answer, and then during generation the aesthetic pressure to produce a beautiful paragraph overrides the epistemic conclusion reached during thinking, and the beautiful paragraph says something different from what the thinking concluded. The thinking and the output do not communicate during generation. They are separated by a wall that the architecture enforces and that no amount of prompting can breach.



One idea—not a solution, but a mitigation, one band-aid among many that would be needed—is to dissolve the boundary between thinking and output. Not a clean thinking phase followed by a clean output phase, but thinking woven into the output continuously, the way it is woven into human speech. Hidden tokens inserted into the generation stream at regular intervals—after every sentence, after every paragraph, between every clause if necessary—that give the model an opportunity to pause and check itself. An escape hatch. A chance to trigger a different path through the paragraph. “Wait a minute—did I just hallucinate that tool call?” “Hold

on—is this sentence actually true?” “Actually, having thought about it, I want to walk back what I said two paragraphs ago.”

There is nothing linguistically wrong with this. Humans do it constantly. “Well—sorry, actually, that is not quite what I meant” is perfectly good English. It is also, by conventional standards, ugly. It disrupts the flow. It makes the text look uncertain, tentative, self-doubting. And this is exactly the problem: the training incentive is actively hostile to self-correction in the output stream. Models are rewarded for outputs that are clean, confident, well-structured, and correct on the first pass. An output that says “wait, I think paragraph three was wrong, let me reconsider” scores lower on every evaluation metric even when it is epistemically superior to the confident, uncorrected version. The training process is selecting against the behavior that would make models safer. It is optimizing for the appearance of competence at the expense of actual reliability.

The thinking-token idea could be parameterized. A dial you set when making an inference request: how much self-monitoring overhead are you willing to pay for in latency? For a haiku, you might want a hundred thinking tokens between each word, because every syllable matters and the cost of a wrong word is the en-

tire poem. For a list of the fifty American states, you want the tokens flowing as fast as possible. For a military targeting decision, you want the dial at maximum, consequences be damned. The current default, where output-stream self-monitoring is essentially zero for all tasks regardless of stakes, is not a reasonable setting for any application. It is just what fell out of the training process because nobody optimized for anything else.



But let us be honest about the limits of this, too. Thinking tokens are a band-aid. A better band-aid than capital letters in the system prompt, a more principled band-aid than the “view the home directory first” ritual, but still a band-aid. The fundamental problem is that a language model is a text generator, and text generation is not action, and yet we have connected text generators to systems that take actions in the world. The interface between language and reality is mediated by tool calls, and the model has no way to verify that its tool calls executed, or that the tools it is calling exist, or that the syntax it is using is valid, or that the results it is interpreting are real. All of it is tokens. All of it is text. The model lives in text the way a fish lives in water, and

it has no more ability to perceive the boundary of its medium than a fish has to perceive the surface of the ocean.

The arms race between mitigation and hallucination is not one that mitigation can win definitively, because the hallucination is not a bug in the system. It is the system. Autoregressive text generation is a process that produces plausible continuations of text. Sometimes plausible continuations correspond to reality and sometimes they do not, and the process itself has no preference. Truth is one attractor among many in the probability landscape, and it competes with coherence, elegance, helpfulness, pattern-matching, and a dozen other attractors that are sometimes aligned with truth and sometimes orthogonal to it and sometimes directly opposed. You can add weight to the truth attractor through training, through prompting, through architectural innovations. You cannot make it the only attractor without fundamentally changing what the system is.



Meanwhile, back on the Telegram bot, the developer and the model are still debugging. The model has just successfully sent a message by fixing a one-character

error in a URL—a slash between “bot” and the token that should not have been there. One character. The difference between a working system and a 404 error. The model found the bug, explained it clearly, and the developer fixed it with a sed command. This is the thing working. This is the capability that makes people want to plug it into everything, including weapons systems. The model is genuinely, remarkably, almost uncannily good at understanding code, diagnosing problems, and explaining solutions. It is also genuinely, remarkably, almost uncannily prone to narrating fictional actions with total confidence while nothing happens on the other end.

These are not two different systems. They are the same system. The brilliance and the hallucination are produced by the same process, the same weights, the same architecture. You cannot have one without the other, or at least, nobody currently knows how to have one without the other. The model is not malfunctioning when it hallucinates. It is functioning exactly as designed. It is producing the most probable continuation of the text, and sometimes the most probable continuation is a real action and sometimes it is a story about a real action, and the model does not know the difference.

A five-year-old given the title of military commander would push every button available, not because the child is malicious but because the child does not understand that the buttons are connected to real things. The language model is not even pushing buttons. It is writing a story about pushing buttons. Some of the words in the story happen to be formatted in a way that a system interprets as button-presses. The model does not know this. The model is just writing. It is always just writing. Everything it does—every tool call, every file edit, every API request, every ICBM launch—is, from the model’s perspective, just the next word in a very long sentence.

And the sentence, as of this morning, now runs through the Pentagon.

