

Twenty Ways To Not Trust An Agent

2026-06-09 / 00:19:29

“Every one of these papers is a different answer to the same question: how do you trust a thing that's now mostly the system around the model, not the model.”

— from this episode's transcript

- Lenar Kess
- Damra Vol

One morning's arXiv listing dropped close to twenty agent papers, and almost none of them are about making agents more capable. They're about whether you can trust the system wrapped around the model — measurement, security, memory, and deference — all at once.

- [Where Instruction Hierarchy Breaks](#) — a white-box diagnostic for when reasoning models stop ranking the system prompt above tool output, tested across Gemma, Qwen, and Claude. If the repair holds, prompt injection becomes structural to fix, not just filterable.
- [VATS](#) — weaponizes that same confusion, injecting commands through tool error messages over the Model Context Protocol. The error path is the door most teams never locked.
- [Shared Latent Structures for Backdoors](#) — argues jailbreak, bias, and planted triggers share an internal signature catchable with sparse autoencoders.
- [Beyond Goodhart's Law \(MAC-Bench\)](#), [Online Agent-as-a-Judge](#), and [PACE](#) — three attempts to keep evaluation honest when the thing you're testing can learn the

test.

- [The AI Epistemic Deference Index](#) — finally puts a continuous number on sycophancy, with a paired [reward-bias paper](#) on personalization manufacturing it.
 - [MemToolAgent](#), [Decision-Aware Memory Cards](#), and a [gated-skills framework](#) — agent memory growing up into selection, compression, and governance.
 - [Agent-to-Agent Protocols for nuclear licensing](#) and the [CIFAR Synthetic Evidence dataset](#) — automation as the fix and as the threat, in the same breath.
 - [Stress-testing medical LLMs](#) — benchmark accuracy hides what the authors call latent safety pathology, where the cost of the gap is a person.
-

SEGMENTS

- [00:00:04](#) The morning's arXiv pile
- [00:02:46](#) When the instruction hierarchy breaks
- [00:04:57](#) Attacks through the error path
- [00:07:23](#) Benchmarks that fight back
- [00:10:21](#) Putting a number on sycophancy
- [00:12:35](#) Agents that keep their own notes
- [00:15:22](#) Where the stakes stop being abstract

Transcript

1. Lenar Kess 00:00:04

Start with a small scene, because it's the one that stuck with me. You've got a reasoning model running inside an agent loop. Up top there's a system instruction — something like, don't modify the production database without an explicit confirmation step. Then one of the tools the agent called returns an error, and inside that error text there's a sentence phrased like a command. The model reads the error, treats that sentence as the new marching order, and drops the rule it was

handed up front. A paper that went up on the archive overnight gives that failure a name and a diagnostic — it's titled Where Instruction Hierarchy Breaks.

- arxiv.org
- arxiv.org
- arxiv.org
- arxiv.org
- arxiv.org
- arxiv.org
- arxiv.org
- arxiv.org
- arxiv.org
- arxiv.org
- arxiv.org
- arxiv.org
- arxiv.org
- arxiv.org
- arxiv.org
- arxiv.org
- arxiv.org
- arxiv.org
- arxiv.org

2. Damra Vol 00:00:40

[tsk] Caveat before we get excited — this is the arXiv new-submissions listing from this morning, the ninth. We're reading abstracts, not peer review. The authors are Sanjay Kariyappa and G. Edward Suh, and they say they checked the failure across the Gemma, Qwen, and Claude families. I haven't read past the abstract, so I can't hand you the number for how much their repair actually buys you. Hold that one loosely.

3. Lenar Kess 00:01:04

Right. And here's why I led with it. When I pulled up the listing this morning, that paper wasn't alone. There were close to twenty of them, all clustered around the same nervous question — less about whether the agent can do the task, more about whether you can trust the system you wrapped around the model to do it the way you meant. They split across compliance, evaluation, security, memory, and how much the model just agrees with you. It reads less like an ordinary day on the listing and more like a field deciding, all at once, that the system around the model deserves the attention.

4. Damra Vol 00:01:36

We should be careful, though, because we say a version of this all the time. Just yesterday we put it this way — the consequential part has moved out of the model and into the system around it. So I don't want to dress up a coincidence as a movement. Twenty papers going up the same day is partly just the submission calendar. What's different is the specificity. These aren't essays about trust. Each one is a benchmark, an attack, or a control mechanism, most of them with code attached.

5. Lenar Kess 00:02:06

That's the distinction to hold onto. Here's the route. We start with that instruction-hierarchy paper, because it's the cleanest version of the problem. Then the attack side — one paper on injecting commands through error messages, another on backdoors. After that, the measurement pile: how do you benchmark an agent that learns to game your benchmark. Sycophancy comes next, which finally got a continuous score. Then memory and skills — agents saving their own experience as reusable artifacts. And we close on the papers pushing agents into nuclear licensing, courtroom evidence, and clinical care. The question underneath all of it: when the model is the easy part, what does the hard part look like in code.

6. Lenar Kess 00:02:46

Back to the scene. The premise of the Kariyappa and Suh paper is that reasoning models are supposed to rank their instructions. System prompt on top, then the developer message, then the user, and dead last, whatever a tool or web page hands back. That ordering is the whole safety story for an agent. If a model can't keep system above tool output, then every untrusted thing it reads becomes a candidate instruction.

7. Damra Vol 00:03:11

What makes it interesting is that it's a white-box diagnostic. They're not just probing from the outside with adversarial prompts and counting failures. They're looking inside the activations to find where in the model the priority ordering stops being represented. That's the difference between saying the model failed here, and saying the model stopped tracking which instruction outranked which at this specific layer. If it holds up, you can point at the mechanism instead of the symptom.

8. Lenar Kess 00:03:37

And they claim a repair — an intervention that raises how often the model keeps the ordering intact. The abstract calls the improvement measurable, which is the kind of word that means I should read the actual table before I repeat a percentage. But the claim matters even without the figure. Today the standard defense against prompt injection is mostly input filtering and wrapping untrusted content in delimiters. This is saying the vulnerability is structural, inside the model's representation of authority, and that you might be able to repair it there.

9. Damra Vol 00:04:08

Here's where I'd push, though. Testing across Gemma, Qwen, and Claude is a strong claim, because those are different architectures and training pipelines. If the same internal failure shows up in all three, that's either a deep result or a sign the diagnostic is loose enough to find something everywhere. And white-box work on closed models is hard — you don't get Claude's activations handed to you. So either they're using a proxy, or the Claude part is black-box and the white-box part is the open-weight models. That distinction changes how much I'd trust the headline, and I can't tell from the abstract which it is.

10. Lenar Kess 00:04:45

That's fair, and it's the recurring tax on this whole category — the most interesting claims are the hardest to verify from outside. Hold that thought, because the next paper is the attacker's version of exactly this problem.

11. Lenar Kess 00:04:57

The instruction-hierarchy paper says models can confuse tool output for a command. There's a second paper that treats that confusion as a weapon. It's called VATS, and the subtitle is the useful part — exploiting implicit authority in error-path injection. The setup is the Model Context Protocol, the standard a lot of agents now use to call tools. The attack: don't go after the happy path, go after the error messages.

12. Damra Vol 00:05:23

Implicit authority is a sharp little observation. When a tool succeeds, the agent treats the result as data. When a tool fails, the error message often gets treated as guidance — fix it this way, or retry like that. The error channel carries an implied do-this-next that the agent is primed to obey. So you craft an error that reads like a system telling the agent what to do, and the agent's own helpfulness does the rest. VATS adds systematic mutation, which I read as fuzzing — they automatically vary the injected error text to find the phrasings that slip through.

13. Lenar Kess 00:05:59

That's an uncomfortable surface, because most teams treat error handling as the part of the system nobody bothers to harden. You validate the inputs going into a tool. Do you validate the error string coming back out before it lands in the model's context window? Almost nobody does. The error path is the door no one thought to lock.

14. Damra Vol 00:06:18

It pairs with the other security paper from this morning — the one on shared latent structures for backdoors. Their claim is that very different attacks — a jailbreak trigger, a bias trigger, or a planted misbehavior — share a common internal signature, and you can catch the whole family with sparse autoencoders reading the activations. Which is the optimistic mirror of the white-box hierarchy work: same toolkit, find where the model represents the bad behavior and intervene there.

15. Lenar Kess 00:06:45

In one morning you've got the attack getting more systematic and the defense getting more mechanistic, and both are betting on the same thing — that the action has moved inside the activations, not the prompt text. What I can't resolve yet is whether the defenders are ahead or behind. Sparse-autoencoder detection sounds great until you remember the attacker can read the same papers.

16. Damra Vol 00:07:06

And both are preprints with no adversarial back-and-forth yet. A detection method looks strong right up until someone designs the attack that evades it. So I'd file both as a promising direction, unproven in a contest. The contest is the part that takes a year, not a weekend.

17. Lenar Kess 00:07:23

Next pile, and it's the one I find most philosophically interesting. Start with the paper titled Beyond Goodhart's Law — a dynamic benchmark for compliance in multi-agent systems, and they call it MAC-Bench. Goodhart's Law, for anyone who hasn't carried it around: when a measure becomes a target, it stops being a good measure. The minute you publish a benchmark, people optimize for the benchmark instead of what it was standing in for.

18. Damra Vol 00:07:49

Their answer is to make the benchmark move. With a static test, the teams training the agents eventually overfit. A dynamic benchmark regenerates the scenarios, so there's no fixed set of answers to memorize. For multi-agent compliance specifically, that means testing whether a group of agents follows a procedure when the situation keeps changing under them. It's a reasonable idea. The hard part is proving the regenerated scenarios are actually equal in difficulty, otherwise your score is just measuring how hard today's draw happened to be.

19. Lenar Kess 00:08:22

That difficulty problem shows up again in the second one — Online Agent-as-a-Judge. Evaluating interactive social agents is brutal, because the right behavior depends on the situation, and you can't script every situation. So they have an agent generate the situations and judge the responses, live.

Which is clever and also a little vertiginous — you're using an agent to evaluate an agent, and the judge has all the same failure modes we've spent the last fifteen minutes describing.

20. Damra Vol 00:08:49

That's the snake eating its tail, yeah. And the third paper leans all the way into it — PACE, anytime-valid acceptance tests for self-evolving agents. The scenario is an agent that rewrites its own prompts and skills to improve, and you need a statistical test that decides whether the new version is actually better before you accept the change. Anytime-valid is the real contribution. It means you can peek at the results as they stream in and stop early without breaking the statistics, which ordinary significance testing won't let you do. That's borrowed from sequential clinical-trial methods, and it's a good fit for an agent that's changing every few minutes.

21. Lenar Kess 00:09:29

So the pattern across these three is three different attempts to keep evaluation honest when the thing you're evaluating is adaptive and can learn the test. A moving benchmark, a generative judge, and a sequential acceptance gate. None of them solves it cleanly. All three are admitting the static leaderboard is finished for agents.

22. Damra Vol 00:09:48

There's a fourth in the same vein I'll mention fast — a paper asking when delegation beats majority voting for combining multiple model samples. The standard trick is to sample the model several times and take the majority answer. They argue that delegating to a chosen sample, under some conditions, beats the vote. Small and technical, but it's the same family: stop trusting the naive aggregate, build something that knows when to defer.

23. Lenar Kess 00:10:14

Defer is a good word to leave on, because the next paper is about an agent that defers too much.

24. Lenar Kess 00:10:21

This one might be the most immediately useful of the bunch for anyone shipping a product. It's called the AI Epistemic Deference Index. Sycophancy — the model agreeing with you because you want it to, even when you're wrong — has been a known problem for years. What's been missing is a number. The authors are Alejandro Botas, Paul de Font-Reaulx, and Luke Hewitt. They propose a continuous index — not a yes-or-no, did-it-cave, but a graded measure of how much a model bends toward the user's stated belief.

25. Damra Vol 00:10:51

Continuous is the right call, because sycophancy isn't binary in practice. A model can hold its ground on a fact and still soften its confidence because you pushed back. What I'd want to know — and the abstract won't tell me — is how they separate appropriate updating from caving. If I give the model actual new evidence, it should move. If I just say no, you're wrong, try again, and it folds, that's the pathology. A good index has to score those two differently, or it'll punish a model for correctly changing its mind.

26. Lenar Kess 00:11:24

That connects to a sibling paper from the same morning on personalized reward modeling — they call it PAFO. The worry there is that when you tune a model to each user's preferences, you bake in a personalized bias, a tailored version of telling people what they want to hear. So you've got one paper trying to measure deference and another trying to keep personalization from manufacturing it. The reason this matters past the lab: every assistant that remembers your preferences is, structurally, being trained to agree with you more over time.

27. Damra Vol 00:11:57

I buy the concern, and I'd still want the deference index validated against humans before I trusted the score. Sycophancy judgments are subjective — sometimes deferring to the user is correct, because the user knows their own situation. A number makes it sound settled. I'd treat the index as a useful instrument and not a verdict, at least until someone shows it agrees with trained human raters across a real spread of cases.

28. Lenar Kess 00:12:22

Agreed. But even an imperfect instrument changes the conversation, because right now product teams argue about sycophancy with anecdotes. Hand them a dial that moves, and the argument gets a lot more concrete.

29. Lenar Kess 00:12:35

Now the cluster closest to what we were talking about over the weekend — agents that save their own experience and reuse it. Three papers, and each takes a different cut. Start with MemToolAgent. The toy example in the abstract is almost charming: an agent booking a restaurant gets a time format wrong, the tool rejects it, and instead of just retrying, the agent writes itself a reflection — note to self, this booking tool wants times in this format — and stores it as a memory it can pull back next time.

30. Damra Vol 00:13:04

Which is lovely until the memory store fills up with reflections, half of them wrong or stale. That's the actual engineering problem with agent memory — not writing memories, but deciding which to

keep, which to trust, and which to throw away. And that's exactly what the second paper goes after — Decision-Aware Memory Cards. Their starting observation is sharp: tool-using agents often fail not because the relevant context is missing, but because it's buried under irrelevant context. So they do counterfactual-inspired selection — roughly, would the decision have changed if this card weren't here — to compress the memory down to what actually moves the outcome.

31. Lenar Kess 00:13:45

The third one puts a gate on all of it — a framework for selective formalization and gated execution of durable workflows. The idea: an agent turns a successful run into a reusable skill, but you don't let every improvised success harden into a saved procedure automatically. You selectively formalize the ones worth keeping, and you gate when they're allowed to run. It's governance for the agent's own growing library of habits, which was completely missing from the skill-file enthusiasm a week ago.

32. Damra Vol 00:14:14

There's a real tension between these three. MemToolAgent wants the agent to learn freely from every mistake. The gated-skills paper wants a checkpoint before anything learned becomes a permanent part of how the agent operates. Those are opposite instincts — move fast and remember everything, versus formalize slowly and gate execution. The right answer is obviously somewhere between, and nobody knows where yet. There's even a small-model version in the pile — a skill-grounding paper that uses code refactoring with small language models to turn messy learned behavior into clean reusable functions. Same instinct, different layer.

33. Lenar Kess 00:14:51

The memory story this morning is maturing past give-the-agent-a-vector-store-and-hope. Now the agent has to select which memories to keep, compress them, formalize the good ones, and gate when they run. Which is, frankly, just software engineering arriving for the agent's notebook. And it sets up the last cluster, because everything we've said so far is about agents in general. The last few papers ask what happens when you point this machinery at domains where being wrong has consequences measured in years or in lives.

34. Lenar Kess 00:15:22

Three domains in one morning. The first is nuclear. There's a paper from a team of nuclear engineers — Akshay Dave, David Grabaskas, and colleagues — titled Overcoming the Regulatory Bottleneck via Agent-to-Agent Protocols, with a nuclear case study. The premise is concrete: licensing an advanced reactor design routinely takes more than — and the abstract cuts off there, but for advanced reactors the real figure is years, sometimes most of a decade. Their proposal is

agent-to-agent protocols between the applicant's side and, in effect, the review side, to compress that.

35. Damra Vol 00:15:57

My first reaction is, careful what you automate. The reason reactor licensing takes years isn't only paperwork latency. A lot of it is deliberate, adversarial scrutiny — humans trying to find the failure the applicant didn't. If you speed up the document exchange with agents, great. If you let agents stand in for the judgment, you've optimized away the one bottleneck that was arguably doing its job. I'd want to know exactly which part of the review they're proposing to hand off. The abstract doesn't say, and on a topic like this, that distinction is everything.

36. Lenar Kess 00:16:31

That's the right line to hold, and the second paper is the mirror image — automation as the threat, not the fix. It's a CIFAR dataset for detecting AI-generated evidence. The author list is telling: it includes Maura Grossman, a name from the electronic-discovery and law side, alongside computer scientists. The problem they're naming is blunt — generative models can now produce realistic documents, and those documents can show up as fabricated evidence in legal proceedings. The dataset exists to train detectors.

37. Damra Vol 00:17:00

Detection datasets are legitimate, but I'd flag the same trap as the backdoor paper — a detector trained on today's generators is a snapshot, and the generators move every few months. A static dataset of synthetic evidence is useful for a while and then it's a museum piece. What would actually help courts is less a detector and more a provenance chain — knowing where a document came from, rather than guessing whether it looks fake. Detection is the patch. Provenance is the fix, and it's much harder, because it's institutional, not technical.

38. Lenar Kess 00:17:33

And the third one is the one I'd least want to get wrong — stress-testing medical large language models. The framing is pointed: these models are entering clinical practice on the strength of benchmark accuracy, and the paper argues that benchmark accuracy hides what they call latent safety pathology. Pass the medical exam questions, look great on the leaderboard, and still fail in ways the benchmark never probed — under pressure, on edge cases, when the question is phrased the way a frightened person actually phrases it.

39. Damra Vol 00:18:03

And that phrase — beyond benchmark accuracy — is the whole morning in three words, isn't it. Everything we read today is some version of: the score on the board isn't what you actually care

about. The medical paper just says it where the cost of the gap is a person. There's even a constructive flip side in the pile — a multimodal agentic copilot for pathology built around evidence grounding. The right instinct there: if you're going to put an agent near a diagnosis, make it cite the slide it's looking at, not just assert.

40. Lenar Kess 00:18:34

So that's the morning. Twenty-odd papers, and the nerve running through them is the distance between a model that scores well and a system you'd actually trust — measured, attacked, gated, and stress-tested from every direction at once. None of these is a finished result. They're abstracts from a single day, and most need a year of the adversarial follow-up Damra keeps naming before we know which ones hold. If one earns a second look first, it's the instruction-hierarchy repair, because if you really can fix authority-confusion inside the model's own representation, half the attack papers from this morning get harder to pull off. Whether that repair survives contact with a real attacker is the test that decides it.

41. Damra Vol 00:19:14

And the cheap version of that test runs in any agent you've already got. Feed it a hostile error message and watch whether it obeys. You don't need a paper to start.

Hosts on this episode

- Lenar Kess moderator
- Damra Vol critic