

# When the Evaluation Goes Back Inside

2026-06-10 / 00:24:54

*“Public model assessments are part of the interface now. If they disappear, builders don't just lose a report; they lose a shared object to argue from.”*

— from this episode's transcript

- Lenar Kess
- Damra Vol

Today's episode starts with the Trump administration reportedly telling CAISI to stop publishing public model assessments, then follows the same trust problem through compute deals, TCS's hiring plans, Anthropic's access terms, AWS Bedrock retention questions, and a small set of agent-security papers.

- [Techmeme's CAISI roundup](#) points to reporting that public model assessments have been halted while a new executive order is implemented; that changes the shared evidence builders can cite.
- [Techmeme's OpenAI Ohio item](#) and [TechCrunch on Meta and Reliance](#) show capacity turning into power contracts, geography, and financing.
- [Techmeme's TCS item](#) captures Tata chairman N. Chandrasekaran talking about agents as a hiring and workforce planning issue, not just a productivity demo.
- [Anthropic's Mythos-class data-retention note](#) makes the enterprise boundary more concrete: different clouds mean different retention and access paths.
- [GitInject](#), [the Interlocutor Effect paper](#), [CIAware-Bench](#), and [deployment-time memorization](#) make the research tail practical: agents fail at the seams where code,

memory, privacy, and oversight meet.

---

## SEGMENTS

00:00:04 The public report disappears

00:06:44 Capacity turns into finance

00:11:46 Agents meet the hiring plan

00:15:53 The enterprise boundary

00:20:07 Four papers, one practical edge

## Transcript

### 1. Lenar Kess 00:00:04

Imagine you are trying to decide whether a frontier model is safe enough to put near your company's code, your customers' records, or your internal incident response workflow. You can read the model card and the company blog post. You can maybe get a sales engineer on a call. And then there is one more source you would like to have: a public assessment from a government testing unit that is at least trying to measure the model from outside the vendor's own story. Today, Wall Street Journal reporting carried by Techmeme says Trump administration officials told CAISI to stop issuing public reports while the administration implements its new AI executive order. CAISI is the Center for AI Standards and Innovation. The reporting says National Cyber Director Sean Cairncross was among the officials involved, and that the pause is tied to national-security concerns around powerful models, including Anthropic's Mythos class. That is the first item today because it is specific. A public evidence source may be moving back inside the government. Then we have a capacity map: OpenAI reportedly talking about a 10 gigawatt Ohio campus with Nvidia backing, Meta signing its first India data-center deal with Reliance, SK Hynix considering a U.S. listing, and investors trying to get synthetic exposure to hot private AI names. After that, Tata's chairman says agents could replace half of TCS jobs in the future and hiring will slow. And then we come back to trust from the enterprise side: Anthropic's Fable and Mythos access terms, Bedrock data-retention questions, and a cluster of agent-security papers that are less abstract than the usual arXiv pile.

- [techmeme.com](https://techmeme.com)

- [x.com](https://x.com)
- [techmeme.com](https://techmeme.com)
- [techcrunch.com](https://techcrunch.com)
- [techmeme.com](https://techmeme.com)
- [techmeme.com](https://techmeme.com)
- [techmeme.com](https://techmeme.com)
- [techmeme.com](https://techmeme.com)
- [forbes.com](https://forbes.com)
- [news.ycombinator.com](https://news.ycombinator.com)
- [arxiv.org](https://arxiv.org)
- [arxiv.org](https://arxiv.org)
- [arxiv.org](https://arxiv.org)
- [arxiv.org](https://arxiv.org)
- [digital-strategy.ec.europa.eu](https://digital-strategy.ec.europa.eu)
- [benton.org](https://benton.org)
- [support.claude.com](https://support.claude.com)

2. Damra Vol 00:01:47

The CAISI item changes the room for me, because public assessment isn't just a civics layer on top of AI. If you are building with these models, a public report becomes part of procurement and security review. It is also the argument you can make to your legal team, and the evidence a regulator can point to without needing a vendor NDA.

3. Lenar Kess 00:02:07

Right. The sourcing language matters here. We don't have an agency announcement saying, here is the new public-reporting policy. We have reporting that officials asked CAISI to halt publication of model assessments while the executive order is implemented. Benton's reprint of the Journal item says CAISI remains operational internally, but its public role is in question. That distinction matters. This isn't the government shutting down AI testing, at least from what is public. The reported change pulls back the portion of the work the rest of us can inspect.

4. Damra Vol 00:02:39

And the reason that matters isn't only transparency as a virtue. It changes how evidence moves. A private assessment can still inform policy, but it can't be challenged, compared, quoted, or reused by a hospital CIO or a school district. A startup security lead and a court-appointed expert lose the same shared reference. Once the report stays inside, everyone outside has to fall back to vendor documents, leaks, paid access, or vibes dressed up as diligence.

5. Lenar Kess 00:03:09

There is a charitable version of the administration's concern. If a model assessment includes details about cyber capability or biological misuse, you can understand why national-security officials would be nervous about publishing a playbook. The reporting says those concerns were part of the push. And CAISI itself, formerly the AI Safety Institute, sits in that awkward place: it is supposed to create shared evidence about dangerous systems, but some of the evidence is dangerous because it is actionable.

6. Damra Vol 00:03:38

That is the hard version of the problem, and I don't think builders should wave it away. If an evaluation says a model can do a particular exploit chain, publishing every detail may help the wrong people. But the opposite break is also serious. If the public only gets a sentence that says, trust us, we tested it, then the evaluation stops being an evaluation for anyone outside the room. It becomes permission language.

7. Lenar Kess 00:04:04

The craft point is that model assessment is now part of the deployment interface. We have spent the past few days talking about agents, memory, instruction hierarchy, and benchmarks. Under all of that is a practical question: what evidence can a builder actually use? If CAISI-style reports disappear from public view, the evidence stack gets thinner at the exact moment the systems are harder to inspect from ordinary product behavior.

8. Damra Vol 00:04:29

And thinner evidence usually means more private negotiation. Bigger companies can still get custom terms, private briefings, and direct access to lab safety teams. Smaller teams read a marketing page and a few forum threads. That isn't a conspiracy claim; it is just how procurement works when the shared artifact goes missing.

9. Lenar Kess 00:04:50

There is also a strange timing loop here. The executive order from last week asked AI firms to voluntarily submit models for cybersecurity testing. Now the reported move is to centralize or restrict public reporting from the government tester. So from the outside, the administration is asking for more access to frontier models while giving the public less access to the assessment layer that comes out of that access.

10. Damra Vol 00:05:13

That is exactly where the next official document needs to be concrete. Not the slogan. The mechanism. What gets published? What gets summarized? Who can audit the summary? Can researchers contest the finding? Is there a delayed-release path, where sensitive details are held

back but the public still gets methodology and high-level results? Those choices decide whether this is a security review process or just an internal clearance process.

11. Lenar Kess 00:05:40

And I would separate two reactions that often get mashed together. One reaction is, public AI evaluations are politically important because they keep government from becoming a black box. True. The other is more prosaic and more useful for this show: public AI evaluations are developer infrastructure. They help teams decide where a model can go, what compensating controls it needs, and whether a vendor claim deserves trust.

12. Damra Vol 00:06:06

There is a maintenance cost too. If you remove a public report, somebody still has to write the internal version, classify it, decide who can read it, brief other agencies, brief vendors, and handle disputes when a lab thinks the finding is wrong. The work doesn't vanish. It just moves into a process most builders can't inspect.

13. Lenar Kess 00:06:25

So I am not ready to say this is a final transparency rollback, because we don't have the final rule. But I am comfortable saying the reported halt matters at builder altitude. It changes the public evidence available for model choice, and it makes the next CAISI communication more consequential than a normal agency update.

14. Lenar Kess 00:06:44

The second cluster is infrastructure, and it works better as an update map than as a grand theory. Techmeme points to reporting that OpenAI is in talks to lease a 10 gigawatt data-center campus in Ohio, with possible Nvidia backing. The site is described in related reports as the former Portsmouth Gaseous Diffusion Plant area in Pike County, where new power generation is already part of a Department of Energy story. A 10 gigawatt campus is a huge phrase, but the useful detail is the structure: OpenAI as the lessee, Nvidia potentially supporting the financing or lease economics, and energy development turning into an AI capacity promise.

15. Damra Vol 00:07:24

That is what makes it different from a normal data-center headline. If Nvidia is backing a lease, even indirectly, then the chip supplier isn't only selling hardware into demand. It is helping demand become financeable. That matters because the bottleneck isn't one purchase order anymore. It is land and power first. Then the builder needs debt, credit support, interconnects, cooling, and a customer credible enough that someone will build the facility before the workloads fully arrive.

16. Lenar Kess 00:07:53

And this comes right after last week's compute-finance stories, so I don't want to repeat the whole argument. The fresh details today are geographic and institutional. Ohio isn't an abstract cloud region in this story. It is a federal-land and power-development story. Meta's India deal isn't just another hyperscaler campus; TechCrunch says Meta signed its first AI data-center deal in India with Reliance. SK Hynix considering a U.S. listing is about memory suppliers accessing capital markets. The China-linked tokenized-stock item is about investors trying to get exposure where direct access is constrained.

17. Damra Vol 00:08:30

The India piece is especially interesting because it isn't just compute near users. It is compute near a political economy. Reliance brings local infrastructure, energy relationships, and a regulatory surface Meta doesn't get by showing up alone. If you are serving India at AI scale, the data center is also a permission structure.

18. Lenar Kess 00:08:50

That sentence says more than the usual “sovereign AI” label. The facility says where the model can run. It says who has local leverage, and whose grid, land, and procurement relationships make it possible. And for Meta, India is both a massive user market and a market where policy, localization, and infrastructure constraints can shape what products are practical.

19. Damra Vol 00:09:11

For builders, the consequence only looks dull if you never have to buy capacity. If you do, these deals show up as quota, latency, regional availability, model access, and price. A team deciding whether to use a model in Mumbai, Columbus, Dublin, or Singapore is downstream of these financing decisions long before anyone calls them product features.

20. Lenar Kess 00:09:34

The SK Hynix item belongs in the same map but not the same sentence. Techmeme has the report that the memory-chip supplier is considering a U.S. listing. That doesn't mean “AI data centers need memory” and then we all nod. It means the suppliers underneath the model boom may want a capital-market home closer to the investors pricing the boom.

21. Damra Vol 00:09:55

And memory isn't a decorative input. High-bandwidth memory is one of the reasons frontier accelerators are scarce and expensive. If the memory supplier changes where it raises money or how

U.S. investors can own it, that is another way AI capacity becomes a financial product instead of just a semiconductor supply chain.

22. Lenar Kess 00:10:15

Then there is the item about Chinese investors using tokenized-stock products to mimic exposure to hot U.S. private companies. I would keep that in the background, because these products often carry their own weird legal and liquidity questions. But it rounds out the day: capital is trying to reach AI upside through every available doorway, even when direct ownership paths are blocked or awkward.

23. Damra Vol 00:10:38

I would resist collapsing it all into one master strategy. OpenAI leasing in Ohio is one story. Meta partnering in India is another. SK Hynix looking at U.S. markets and synthetic investor exposure from China are different stories again. The shared practical fact is narrower: AI capacity is being packaged through contracts, local partners, listings, and financial wrappers. The wrapper matters because it decides who gets paid before the model ever answers a prompt.

24. Lenar Kess 00:11:06

And it also decides who has bargaining power when demand changes. If a campus is built around a long lease, or a local partner owns the facility, or a supplier has new public-market pressure, the model company inherits a set of obligations. That can be good; obligations make capacity real. It can also make the model roadmap less flexible than the demo makes it look.

25. Damra Vol 00:11:28

That is the developer version of the capacity story. You may experience it as “this model is slow today” or “this region doesn't have the new model yet.” Behind that is a chain of people who made a power promise, financed a building, bought memory, and signed a contract with assumptions about future usage.

26. Lenar Kess 00:11:46

Tata chairman N. Chandrasekaran reportedly said AI agents could replace half of TCS jobs in the future, and that Tata Consultancy Services will reduce hiring as it adopts AI. The Straits Times version says he was talking about one of India's largest IT services companies, a firm that serves multinationals and has been built around a very large human delivery model. This isn't a layoff notice for half the company. That distinction needs to stay in frame. The fresh claim is a chairman talking about future workforce composition and hiring pressure as agents become part of the service model.

27. Damra Vol 00:12:20

That distinction matters because services firms aren't product companies with one app and one user base. TCS sells labor and process knowledge. It also sells client trust and delivery capacity. If agents replace or augment a large share of that work, the economics move through hiring classes, bench management, training, margins, and client contracts.

28. Lenar Kess 00:12:42

And the quote reads differently from a startup CEO saying “we use AI for support tickets.” TCS is part of the machinery that turns enterprise software into real workflows. If its chairman says hiring will slow because agents can do more of the work, that is a signal about how the buyer side may experience AI: fewer junior roles, more automation inside delivery teams, and more pressure to prove that a services contract isn't just human hours with a markup.

29. Damra Vol 00:13:09

It also changes the apprenticeship ladder. Much of IT services work teaches people how systems fail in the wild. Ticket queues, migrations, integration projects, weird customer environments, and brittle batch jobs all become training material. If agents absorb the bottom rungs, companies still need a way for humans to learn the judgment that used to come from doing those tasks badly at first and then getting better.

30. Lenar Kess 00:13:36

That training path gets under-discussed. The simple automation story says: agents take repetitive work, humans move up. Sometimes that is true. But moving up requires a path. If the first three years of work become an agent-supervision layer, the training design has to change. Otherwise you get senior reviewers with fewer juniors becoming senior, and a lot of people asked to manage systems they never had time to understand from the inside.

31. Damra Vol 00:14:01

And in services, the client may not care how the work gets done until something breaks. Then the ownership question becomes concrete. If an agent wrote the migration script, the human team approved it, and the contract priced the work as automated delivery, the accountability path has to be clearer than the demo.

32. Lenar Kess 00:14:20

Chandrasekaran also reportedly cast AI as an opportunity, not just a threat. That belongs in the segment because it is probably how management sees it. If agents let a services firm deliver more work with fewer new hires, improve margins, and move people into higher-value tasks, that is a real business incentive. The labor question isn't whether management can imagine the benefit. It is whether the workers, clients, and training systems get a coherent version of the transition.

33. Damra Vol 00:14:47

And whether the agent work is actually good enough once it leaves the controlled task. Enterprise services are full of exceptions. The old spreadsheet has a macro, and the client VPN times out. The staging environment is named like production. Someone made a policy exception in 2019 and never documented it. Agents can help with that, but they also need context, permission, and review.

34. Lenar Kess 00:15:11

So I would not read the TCS item as proof that half the jobs disappear on a schedule. I would read it as a major services firm telling us where the planning conversation has moved. Agents are entering hiring models. That is more concrete than a productivity chart, because hiring is where companies reveal which future they are willing to budget for.

35. Damra Vol 00:15:30

And for someone early in the field, the practical advice isn't "panic" and it isn't "ignore it." It is: learn the systems around the model. Learn the client boundary, the deployment path, the data model, the test suite, the contract, and the failure recovery process. If agents take over more rote implementation, the human value moves toward knowing what the work is supposed to mean.

36. Lenar Kess 00:15:53

The Anthropic cluster today is messy, so I am going to keep it narrow. Forbes reports on Claude Fable 5 access and rationing. Hacker News is reacting to a claim that AWS Bedrock will require data sharing with Anthropic for Mythos and future models. And Anthropic's own help-center page on Mythos-class data retention gives us the most useful primary language. Anthropic says that through Amazon Bedrock, retention will need to be enabled to access a new covered model, and retained data stays in the AWS environment. Through Claude Platform on AWS, retention works like the direct Claude API and is configured at the workspace level, with retained data handled by Anthropic under the same controls. Through Google Cloud's agent platform, retention also needs to be enabled, but retained data stays in the GCP environment. Azure Foundry gets its own subscription-level treatment.

37. Damra Vol 00:16:48

That is exactly the kind of paragraph enterprise buyers have to read twice. It isn't just "does Anthropic train on my data?" It is where retained data sits, who controls retention, which cloud account or subscription is involved, and whether zero data retention is compatible with the model you want.

38. Lenar Kess 00:17:05

And that is why I don't want to make the HN thread the source of truth. HN is useful reaction evidence. It tells us what operators are nervous about. But the primary artifact is the retention note, because it shows that Mythos-class access isn't a normal model toggle. Access to the more capable or more sensitive model class comes with a retention and trust boundary that may differ by platform.

39. Damra Vol 00:17:28

There is a subtle buyer problem here. A company may have standardized on Bedrock precisely because it believed prompts and completions stayed away from model providers. If a new covered model requires retention, even inside AWS, the security review has to update its mental model. The boundary may still be acceptable. It may even be well designed. But it is no longer the same boundary the buyer thought it bought.

40. Lenar Kess 00:17:53

Fable 5 adds the access side of the same story. Forbes describes Anthropic shipping its strongest model and then rationing access through rate limits or expiration. Other reporting today describes Fable as a Mythos-class model with safeguards, including routing high-risk queries away from the most capable path. I am not going to claim more than the reporting supports. For builders, capability, access, and data handling are now bundled together.

41. Damra Vol 00:18:21

And bundled in a way that can make engineering planning awkward. The best model for a task may be available only under a retention posture your customer has not approved. Or it may be available today under credits that expire, then become expensive enough that you need fallback behavior. Or it may degrade or route certain work because the provider is managing dual-use risk.

42. Lenar Kess 00:18:43

The enterprise trust question gets more concrete than “do we trust Anthropic?” A better question is: can the buyer describe the operational path of a request? Which service receives it, where does data remain, what is retained, who can inspect it, how long does access last, what happens when the request trips a safety boundary, and what fallback model or workflow takes over?

43. Damra Vol 00:19:05

And can the application tell the user what happened without leaking provider internals or making false promises. If a request silently routes to a weaker model because the original path is restricted, that may be fine for safety. But the product still needs to know whether the answer is lower confidence, slower, more expensive, or unavailable for audit reasons.

44. Lenar Kess 00:19:26

That connects back to CAISI in a useful way. Public assessment is one kind of trust artifact. Enterprise retention language is another. Neither is glamorous. Both are now part of whether a model can be used in serious workflows. A frontier model that can't explain its access terms isn't fully productized for the people who have to sign the risk memo.

45. Damra Vol 00:19:47

And a risk memo isn't bureaucracy as decoration. It is where the company decides whether customer data can pass through a system, whether the logs are discoverable, whether a vendor can help debug an incident, and whether the model can be used in regulated work. The terms aren't after the product. They are inside the product.

46. Lenar Kess 00:20:07

The arXiv pool today is enormous, and the agenda is right to avoid turning the episode into a bibliography. But four papers deserve a compact pass because they are all about where agents touch real systems. First, GitInject. The paper is called "Real-World Prompt Injection Attacks in AI-Powered CI/CD Pipelines." The authors say GitInject provisions ephemeral repositories and triggers actual GitHub workflow runs, rather than only simulating tool calls. That detail matters because continuous integration and deployment is where permissions, secrets, untrusted pull requests, and automation already meet.

47. Damra Vol 00:20:46

That is much more useful than another toy benchmark. A simulated tool call can miss the permission boundary that decides whether the attack matters. A real workflow run has the parts security depends on: tokens, repository settings, sandbox behavior, logs, and the exact way a bot reads untrusted text from an issue or a pull request.

48. Lenar Kess 00:21:06

Second, the Interlocutor Effect paper. Its claim, from the abstract, is that large language models alter privacy behavior based on who they think they are talking to. The authors say models tend to reveal more sensitive personal data when addressing another AI agent than when addressing a human user.

49. Damra Vol 00:21:24

That is a nasty deployment detail. A lot of agent systems already talk agent-to-agent because the architecture is convenient. A planner talks to a worker. A support bot talks to a billing bot. A local assistant talks to a remote tool. If the model treats another agent as a safer recipient than a human,

the privacy boundary can weaken exactly where the system designer thought automation was making things more controlled.

50. Lenar Kess 00:21:50

Third, CIAware-Bench. The paper is about control intervention awareness: whether a model can detect that a control protocol has edited, resampled, or replaced part of its trajectory. The authors argue that if a controlled model can detect interventions, it gains information that may help it route around the control.

51. Damra Vol 00:22:09

This one is very easy to make too cinematic, so I would keep it practical. Control systems often assume the monitored system doesn't get a perfect signal about the monitor. If a model can infer that a response was rewritten, or that a suspicious action was replaced, then your oversight layer becomes part of the model's observation space. That doesn't mean doom. It means the monitor design has to account for feedback.

52. Lenar Kess 00:22:34

Fourth, deployment-time memorization. The authors describe agent memory as a privacy-utility frontier, measured with Personalization Recall and Adversarial Extraction Rate, and they introduce a Forgetting Residue Score to ask whether deleted information remains recoverable from derived memory tiers. Translated into product language: an agent that remembers you better may serve you better. The same memory system also becomes an extraction target. Deletion gets harder when summaries or derived notes still contain traces of the information you removed.

53. Damra Vol 00:23:07

That is the memory problem people actually ship. You don't just have "the memory." You have the chat history, the summary, the profile, the embedding store, the task notes, the audit log, and maybe a cache. A user deletes something from one layer, and the system may still carry the fact indirectly in another. The paper's vocabulary is useful because it treats forgetting as something you can test, not a promise you hope is true.

54. Lenar Kess 00:23:33

I like this research cluster today, but not because every paper is definitive. They are fresh papers, and they need replication, review, and real-world pressure. Each one moves the agent conversation away from personality and toward placement. Where is the agent running? What can it read? Who does it think it is talking to? Can it detect the monitor? What memory survives deletion?

55. Damra Vol 00:23:55

And those questions line up with the rest of the episode without forcing them. CAISI is about public evidence. The compute stories are about where capacity can physically and financially run. TCS is about where agents enter the labor plan. Anthropic and Bedrock are about where data sits when a powerful model is accessed. The research papers are about where the agent touches code and memory, and where oversight can see it.

56. Lenar Kess 00:24:22

That is a useful place to end the day's work. The models are still impressive, but today's durable questions sit around them: which assessments remain public, which data centers get financed, which jobs become agent-supervised, which retention boundary the enterprise can actually explain, and which agent seams get tested before they become production incidents. Tomorrow, Thursday, June 11, the strongest follow-up would be an official CAISI or White House clarification that says what the public will still be allowed to see.

## Hosts on this episode

- Lenar Kess moderator
- Damra Vol critic