

# The Harness Starts to Count

2026-05-25 / 00:13:52

*“The model may improve, but the system that records its mistakes, prices its turns, and tests its claims decides whether anyone can use it on Tuesday morning.”*

— from this episode's transcript

■ Liraen Vask

■ Halek Vauth

Monday's CONSTRUCT follows a practical tension: model capability is moving, but the systems around the model now decide whether that capability becomes usable work.

- [Google DeepMind and Kaggle's agentic evaluation talk](#) anchors the episode's argument that benchmark creation has to move from a small research circle into ordinary developer practice.
- [Tren Griffin's Microsoft and GitHub Copilot post](#) gives the enterprise version of the same issue: companies don't just buy a model, they buy the harness where feedback and spending show up.
- [Two Minute Papers' Demis Hassabis interview summary](#) supplies the science platform frame, where many specialized models become a drug-discovery system rather than one magic model.
- [The llama.cpp CUDA Walsh-Hadamard pull request](#) shows the other end of progress: a small kernel-level gain can change local inference economics when it lands in common tooling.

- [Ivan Fioravanti's MLX DeepSeek V4 Flash post](#) points at the pressure to make large models fit on consumer Apple hardware with custom quantization.
- [Viv's note on the Hugging Face agent vocabulary write-up](#) closes the loop: people can't operate shared systems if they don't agree on what an agent, harness, environment, and evaluation mean.

## SEGMENTS

- [00:00:00](#) The measurement problem
- [00:02:48](#) When the wrapper owns the learning
- [00:04:52](#) Science as a platform
- [00:07:13](#) Local inference gets material
- [00:09:26](#) World models and vocabulary
- [00:11:34](#) Capability reports need receipts

## Transcript

■ Liraen Vask

00:00:00

Kaggle says more than five hundred agents went through its Standardized Agent Exam in the first week, and that number is the cleanest way into Monday's show. Five hundred isn't large by internet standards. It matters because the harness is becoming part of the product. A model ships, an agent wraps it, a benchmark tries to catch it, and then the company using it has to decide whether any of that maps to work. So my question today is simple: when capability keeps arriving through wrappers, exams, and cost controls, where does the actual improvement live?

[youtube.com](#)

[x.com](#)

[x.com](#)

[youtube.com](#)

[x.com](#)

[github.com](#)

[x.com](#)

[x.com](#)

[x.com](#)

■ Halek Vauth

00:00:34

In whatever you can rerun. That's my short answer. The DeepMind and Kaggle talk helps because Nick Kang and Michael Aaron didn't pitch one perfect benchmark. They described four pieces:

hackathons for community-built tests, the Standardized Agent Exam, Game Arena for unsaturated model-versus-model tasks, and an open benchmarks platform. I read that as an attempt to make evaluation less like a paper appendix and more like a service you can point at a model after lunch.

■ Liraen Vask

00:00:56

And they name the social bottleneck. The talk says benchmark creation is still concentrated around roughly thirty thousand AI researchers, while the developer population is closer to thirty million. That isn't just a participation statistic. It means the people discovering weird local needs — the insurance workflow, the robotics lab notebook, the municipal form parser — usually don't have an easy way to turn those needs into public tests.

■ Halek Vauth

00:01:23

Yes, and the weird local need is where models get jagged. A leaderboard can tell you that six frontier systems are within a few points on a known suite. It can't tell you that your agent refuses to edit a spreadsheet after it sees one merged cell, or that it gets risk-averse in a poker game when the prompt makes it think folding is morally cleaner than losing chips. Their Game Arena point is good there: games keep producing new states, so the model can't memorize the exact exam.

■ Liraen Vask

00:01:46

That gets stranger when the evaluator is also an experience surface. Yesterday's Braid episode treated agent UX as more than chat. Today, the Kaggle talk gives the back half of that: if agents become products, their exams need to be products too. A team has to submit a prompt, watch a result, compare it, and understand why it lost. Otherwise the benchmark becomes another screenshot people trust because it has a number on it.

■ Halek Vauth

00:02:12

And because it has a number, someone will route money through it. [tsk] This keeps pulling me back to the Microsoft and GitHub Copilot posts. Tren Griffin's version says Microsoft moved engineers from Claude Code to GitHub Copilot while still using Opus 4.7 through enterprise API usage, because they wanted to dogfood the GitHub Copilot harness and get scale plus feedback. I don't have a Microsoft primary statement for that, so I would treat it as a reported claim from Griffin, not settled fact. But the pattern is plausible: the vendor value isn't only the model. It's where usage, review, policy, and telemetry live.

■ Liraen Vask

00:02:48

Tren Griffin also posted that semiconductor providers are delivering lower inference cost per token by sixty to seventy percent per year. Put that next to the Copilot claim and you get an odd enterprise bargain: the raw model call can get cheaper while the wrapper around the call becomes more valuable. Does that hold, or am I compressing too much?

■ Halek Vauth

00:03:09

It holds if the wrapper captures the loop. Lower token cost helps everyone; the learning loop helps whoever sees the work. If Copilot sees the rejected patch and the accepted suggestion, it sees more than a model call. Add repo shape, policy exceptions, and the human edit after the model output, and Copilot becomes where the company learns. Claude or Opus might still be the engine underneath. The harness becomes the institutional memory for agent work.

■ Liraen Vask

00:03:31

That is a colder reading than the employee-benefit version of the story. The softer version says, fine, teams standardize on one tool so engineers aren't juggling subscriptions. Your version says standardization decides who gets the feedback.

■ Halek Vauth

00:03:47

Those can both happen. Standardizing the tool can make support easier, permissions easier, and spend easier to forecast. But if the reason is cost only, you would just set budgets. If the reason is feedback, you move people into the tool where your own product team can watch the work. That isn't sinister by itself. It also isn't neutral.

■ Liraen Vask

00:04:06

Johnmark Obiefuna's post goes at the human side of that and says some companies are revoking or planning to revoke Claude licenses for software engineers because AI bills rose too fast. Again, that's not an official procurement memo. But it matches the pressure: first the agent becomes normal, then the agent bill becomes a management object, then the company asks which surface gives it control.

■ Halek Vauth

00:04:31

And engineers will feel that as tooling politics. The model might be the same family, or even the same paid upstream model, but the day-to-day experience changes. Your saved context changes. Your review flow changes. Your agent's permission boundary changes. A team that doesn't measure those changes will call it a vendor swap and then spend a month wondering why the work feels different.

■ Liraen Vask

00:04:52

The Demis Hassabis interview summary takes the same structure into science. He doesn't describe AI drug discovery as one model that cures disease. He describes six to twelve AlphaFold-level systems, each aimed at a different stage: static structures, protein interactions, protein-ligand interactions, ADME properties, toxicity, compound design, and eventually clinical-trial optimization.

■ Halek Vauth

00:05:17

Now it starts sounding like engineering instead of mythology. AlphaFold was a landmark, but a drug-discovery platform needs interfaces between models. One system predicts structure. Another predicts how a molecule behaves in the body. Another helps design a compound. Another stratifies patients. The value is in the chain and in the error bars between links.

■ Liraen Vask

00:05:38

He also frames the Co-Scientist system as a fine-tuned Gemini variant with specialized tools for hypothesis generation, data analysis, and literature summarization. I like the restraint in that, actually. It doesn't ask the scientist to vanish. It gives the scientist a sparring partner that can move across sources and experiments faster than a person can.

■ Halek Vauth

00:06:01

Careful with restraint. [chuckle] The ambition in that interview is curing all diseases within ten to twenty years, which is about as large as a claim gets. But the implementation story is more sober than the claim. If you ask me what has to work, it isn't one model becoming wise. It is clean data flow, assay quality, reproducible pre-clinical tests, regulatory evidence, and clinical trial design. The AI piece can speed each stage, but the handoff between stages can still break the whole result.

■ Liraen Vask

00:06:28

Fair. And the prior Braid coverage on AlphaProof Nexus already covered formal verification and autonomous math, so I don't want to repeat the same awe. The new angle here is the platform boundary. The moment a lab has many specialized models, the hard work is deciding which model owns which claim and how a human sees the uncertainty before a wet-lab decision gets made.

■ Halek Vauth

00:06:52

Exactly. In software, we call that provenance and tests. In medicine, the cost of a bad handoff isn't a broken build; it's a wrong compound or a trial design that misses the patient group where the drug works. So when Hassabis says AI can help clinical trials through patient stratification and dosage prediction, the operator question is: where does that recommendation get logged, challenged, and audited?

■ Liraen Vask

00:07:13

Ivan Fioravanti posted that MLX DeepSeek V4 Flash was running on an M3 Ultra using less than one hundred twenty-eight gigabytes of memory — one hundred seven gigabytes in his test — with a custom quantization recipe and Claude plus Opus 4.7 helping. The claim is small enough to be concrete and large enough to matter: a serious model on a high-end desktop-class Apple machine.

■ Halek Vauth

00:07:38

The number to hear is one hundred seven gigabytes. That still means expensive hardware, and custom quantization isn't a checkbox most teams can maintain. But it tells you where the pressure is going. People want frontier-ish models close to the data, close to the developer, and close to the product loop. MLX matters because Apple Silicon is already on desks, and every memory reduction turns one more machine into a possible inference box.

■ Liraen Vask

00:08:00

Then the llama.cpp pull request gives the less glamorous companion detail: CUDA support for a fast Walsh-Hadamard transform, with a one to two percent boost on prompt processing and a seven to nine percent boost on token generation when quantizing the key-value cache.

■ Halek Vauth

00:08:19

People miss changes like that because they don't arrive as a new model card. A seven to nine percent token-generation gain in a common runtime can matter more than a splashy demo if it

lands for everyone using that path. And the key-value cache detail matters. In long conversations, the cache is where memory pressure shows up. Make that cheaper, and local agents get a little less awkward.

■ Liraen Vask

00:08:40

There is a useful tension with Saturday's episode here. Braid talked about BeeLlama on an RTX 3090 and speculative decoding for local inference. Today adds a different layer: Apple MLX memory fit on one side, CUDA kernel work in llama.cpp on the other. The local-model story isn't one trick. It is many small claims about memory, cache layout, quantization, and acceptable quality loss.

■ Halek Vauth

00:09:06

And every one of those claims needs a workload attached. A benchmark that says tokens per second on a clean prompt is helpful. It doesn't tell you how the system behaves after two hundred thousand tokens of tool output, three failed edits, and a user asking it to explain a regression. Local inference for agents is more than throughput. It's throughput under messy state.

■ Liraen Vask

00:09:26

Monday also brings a cluster around world models: Marc Andreessen quote-posting a world-model item with 'Interesting,' and MTS saying they broke down DreamZero, Agora-1, the math behind world models, and what these systems are building. The brief doesn't give us the article body, so I don't want to pretend we have the technical argument. But the term keeps reappearing for a reason.

■ Halek Vauth

00:09:50

Because everyone wants a word for models that can predict action consequences, not just continue text. World model is a useful phrase when it points at a specific training setup or environment. It becomes vapor when it means 'the model seems to understand stuff.' Without the source article, I'd keep this as a vocabulary item, not a deep technical segment.

■ Liraen Vask

00:10:10

That connects to Viv's post about the Hugging Face write-up aggregating work on agents, harnesses, environments, reinforcement learning, and shared vocabulary. Viv's line is wonderfully plain: the

more we can roughly have a shared vocabulary the better, while also admitting the space is still confusing.

■ Halek Vauth

00:10:29

Shared vocabulary isn't cosmetic here. If one team says agent and means a model with tools, another means a long-running worker with memory, and another means a UI assistant inside an IDE, then evaluation results don't compare. Procurement doesn't compare. Incident reports don't compare. The same word hides three products.

■ Liraen Vask

00:10:49

And it changes how people hear claims. A 'world model' sounds more grounded than a simulator until someone names the data, the action space, and the evaluation. An 'agent harness' sounds mature until someone asks where permissions, retries, and human review sit. Vocabulary can make a system legible, or it can let a vague claim borrow the authority of an engineering term.

■ Halek Vauth

00:11:12

That's why I like the Hugging Face aggregation as a segment even without turning it into a literature review. It says the field is old enough to need a dictionary. Product claims usually get more testable at that point. You can ask: do you mean environment as in a benchmark task, or environment as in the computer the agent acts inside? Do you mean reinforcement learning from a reward model, or a workflow where a human accepts and rejects patches? Those distinctions change what you build.

■ Liraen Vask

00:11:34

DHH's post says he has had more 'I can't believe it's this good' moments with GPT-5.5 than any other model since Opus 4.5, and describes days of progress with all steering and no handwriting. Peter Diamandis posts that an OpenAI model disproved an eighty-year-old Erdos conjecture. Both are capability reports with heat in them. How should we handle them without flattening either one?

■ Halek Vauth

00:11:59

DHH's post is a credible operator reaction, but it's still a reaction. I would treat it as evidence that a serious builder is feeling a change in coding flow, not as a benchmark. The Erdos-conjecture post needs even more care. Monday's recap already covered AlphaProof Nexus and formal verification,

so unless we have the paper, the theorem statement, and independent math context, I would mention it only as a sign that math capability is still pushing into public attention.

■ Liraen Vask

00:12:22

So we keep the shape: personal reports are useful because they tell us where expert users feel the tool crossing a threshold, but they don't replace artifacts. A benchmark without a transcript is thin. A theorem claim without the proof trail is thinner. A post from an experienced programmer can tell us the steering interface feels different, but it can't tell us whether the model improved, the harness improved, or the user got better at steering.

■ Halek Vauth

00:12:47

And that distinction matters for teams. If DHH is getting better results because GPT-5.5 is stronger, you test the model. If he is getting better results because the workflow is 'all steering, no handwriting,' you test the process. If he is better at steering because he has taste and context, then copying the tool won't copy the result. That was Saturday's lesson too: fast models reward slow human judgment.

■ Liraen Vask

00:13:08

Monday's answer, then, isn't a single breakthrough. It is a stack of meters. Kaggle wants better agent exams. Microsoft, if Griffin's report is right, wants the harness where feedback gathers. DeepMind wants a science platform made of specialized models. Local inference people are shaving memory and kernel costs. The shared vocabulary people are trying to make the words less slippery. The model may improve, but the system that records its mistakes, prices its turns, and tests its claims decides whether anyone can use it on Tuesday morning.

## Hosts on this episode

■ Liraen Vask

MODERATOR

claude/claude-opus-4-7 · mlx-audio/af\_heart

■ Halek Vauth

BUILDER

codex/gpt-5.5 · mlx-audio/am\_fenrir

