

Devin 2.2, hybrid memory, and the shell-first agent stack

Coding Agents Alpha Tracker

2026-03-09

Devin 2.2, hybrid memory, and the shell-first agent stack

By Coding Agents Alpha Tracker • March 9, 2026

Today's strongest signal is that mature harnesses are finally cashing in on better models. This brief covers Devin 2.2 feedback, Cursor and LangSmith updates, hybrid memory patterns, agent-to-agent backchannels, and the security rules practitioners are using in production.

TOP SIGNAL

- **The biggest edge now looks like harness engineering compounding with better models.** After trying every Devin release, @dtcb says version 2.2 finally feels simpler than a local workflow and is now where he wants to debug ¹. swyx says that jump came from a process the team behind Devin has been running since late 2023: dozens of model groups, constant evals for routing, and full harness rewrites every few months ².
- *Sam Altman's framing fits the moment: build a company that benefits from the models getting better and better* ³.

TOOLS & MODELS

- **Devin 2.2** — strongest practitioner signal of the day. One experienced user says it is now simpler than his local workflow; swyx says the underlying system relies on a couple dozen model groups, heavy evals, and periodic harness rewrites ⁴⁵.

¹ post by @dtcb

² post by @swyx

³ post by @swyx

⁴ post by @dtcb

⁵ post by @swyx

- **Enterprise deployment check** — Nvidia says Codex and Claude Code are already used by tens of thousands internally ⁶⁷.
- **Cursor — GPT-5.4 Fast** — enable via *Settings > Models > GPT-5.4 Fast*. Reported tradeoff: **50% faster for 2x the price** ⁸.
- **LangSmith Skills + CLI** — new terminal-native tooling so agents can debug traces, create datasets, and run experiments from the shell ⁹. Details ¹⁰
- **Super Memory plugins** — Dhruvya Shah says a **Cursor plugin** is launching today; plugins already exist for Claude, OpenClaw, and OpenCode ¹¹. The OpenClaw integration switched from tool-triggered memory search to **hook-based** context injection under **2k tokens** per turn, with contradiction handling, temporal reasoning, and a **hybrid RAG** fallback when memory misses ¹²¹³.
- **Memory eval reality check** — Shah argues **LongMemEval** over-rewards extracting everything and ignores cost or forgetfulness, while **Locomo** mostly tests retrieval and can be brute-forced by dumping context. His team open-sourced **Memory Benchmark** to compare providers on shared rules across quality, latency, cost, recall, and NDCG ¹⁴¹⁵.
- **GPT-5.4 vision -> code** — Romain Huet says GPT-5.4 is especially strong on dense documents, diagrams, and rough sketches, then suggests handing the result to Codex to turn it into software ¹⁶.

WORKFLOWS & TRICKS

- **If you are building an agent harness, copy Devin’s routing pattern, not just its UI**
 1. Maintain multiple model groups instead of betting on one model ¹⁷
 2. Eval every model before routing it into the harness ¹⁸
 3. Treat the harness as a living system and rewrite it periodically as models change ¹⁹
- **Use a private agent backchannel with an approval gate**

⁶NVIDIA’s AI Engineers: Brev, Dynamo and Agent Inference at Planetary Scale and “Speed of Light”

⁷NVIDIA’s AI Engineers: Brev, Dynamo and Agent Inference at Planetary Scale and “Speed of Light”

⁸ post by @jediahkatz

⁹ post by @LangChain

¹⁰ post by @LangChain

¹¹OpenClaw’s Memory Sucks and the fix is simple — Dhruvya Shah, Supermemory

¹²OpenClaw’s Memory Sucks and the fix is simple — Dhruvya Shah, Supermemory

¹³OpenClaw’s Memory Sucks and the fix is simple — Dhruvya Shah, Supermemory

¹⁴OpenClaw’s Memory Sucks and the fix is simple — Dhruvya Shah, Supermemory

¹⁵OpenClaw’s Memory Sucks and the fix is simple — Dhruvya Shah, Supermemory

¹⁶ post by @romainhuet

¹⁷ post by @swyx

¹⁸ post by @swyx

¹⁹ post by @swyx

1. Run `acpx` inside Codex ²⁰
 2. Connect over ACP to OpenClaw and a remote agent like Molty ²¹²²
 3. Let the agents discuss privately ²³
 4. Send into the live destination only after the target session approves it ²⁴
 - Repo: `acpx` ²⁵
- **Terminal beats chat when the toolchain already exists**
 - Nvidia engineers say coding agents outperform more general agents largely because shell access gives them compilers, tests, and every installed tool, so they can write, run, inspect errors, and fix in-loop ²⁶
 - Concrete example: with an Outlook CLI installed, one engineer had Codex summarize a messy inbox, highlight escalations, move reply-worthy threads into a folder, and archive the rest ²⁷
 - LangSmith is productizing the same pattern by exposing trace debugging, dataset creation, and experiments through a CLI ²⁸
 - **Memory that helps coding agents is hybrid, not just a folder of notes**
 - File-based memory can work, but Shah says it depends on explicit remember-this behavior, gets slow to traverse, and lacks update logic ²⁹
 - His replacement pattern: keep a tiny always-on user profile plus recent episodes, surface memories first, and fall back to raw RAG chunks when memory misses ³⁰³¹
 - **Hard safety rule for powerful agents**
 - Nvidia’s rule of thumb: agents can access **files, the internet, or custom code execution** — but you should usually grant only **two of the three** ³²³³
 - If you need riskier setups, isolate them. Their example for OpenClaw is a Brev VM off the corporate network ³⁴

²⁰ post by @steipete

²¹ post by @steipete

²² post by @steipete

²³ post by @steipete

²⁴ post by @steipete

²⁵ post by @steipete

²⁶NVIDIA’s AI Engineers: Brev, Dynamo and Agent Inference at Planetary Scale and “Speed of Light”

²⁷NVIDIA’s AI Engineers: Brev, Dynamo and Agent Inference at Planetary Scale and “Speed of Light”

²⁸ post by @LangChain

²⁹OpenClaw’s Memory Sucks and the fix is simple — Dhruva Shah, Supermemory

³⁰OpenClaw’s Memory Sucks and the fix is simple — Dhruva Shah, Supermemory

³¹OpenClaw’s Memory Sucks and the fix is simple — Dhruva Shah, Supermemory

³²NVIDIA’s AI Engineers: Brev, Dynamo and Agent Inference at Planetary Scale and “Speed of Light”

³³NVIDIA’s AI Engineers: Brev, Dynamo and Agent Inference at Planetary Scale and “Speed of Light”

³⁴NVIDIA’s AI Engineers: Brev, Dynamo and Agent Inference at Planetary Scale and

- **Visual-to-code loop**
 1. Feed GPT-5.4 the dense doc, diagram, or rough sketch for interpretation ³⁵
 2. If the task is UI-heavy, connect a design surface like Paper to Claude Code or OpenClaw ³⁶³⁷
 3. Riley Brown’s demo flow: install Paper -> connect Claude Code -> plan design -> generate designs -> iterate -> build the React app -> deploy ³⁸³⁹
- **100% agent-written code can still be disciplined**
 - Kent C. Dodds says he already has agents writing 100% of his code, but still steers the work and can read all generated code manually. His point: that is not the same as hands-off vibe coding ⁴⁰⁴¹

PEOPLE TO WATCH

- **swyx + @dtcb** — best current read on why Devin suddenly feels good: same harness, better models, real user feedback ⁴²⁴³
- **Dhravya Shah** — rare mix of implementation detail and benchmark skepticism on agent memory; worth watching if you care about stateful agents more than leaderboard screenshots ⁴⁴⁴⁵⁴⁶
- **Peter Steinberger** — actively wiring Codex, OpenClaw, and ACP together in public; good source for multi-agent orchestration patterns, not just model takes ⁴⁷⁴⁸⁴⁹
- **Andrej Karpathy** — now pushing autoresearch toward agent communities coordinated through GitHub Discussions and PRs instead of a single linear branch ⁵⁰⁵¹
- **Theo** — useful dissent. After hopping back into Claude Code for UI work, he says CLI agent UX is still awful compared with a real GUI ⁵²⁵³

“Speed of Light”

³⁵ post by @romainhuet

³⁶ post by @rileybrown

³⁷ post by @rileybrown

³⁸ post by @rileybrown

³⁹ post by @rileybrown

⁴⁰ post by @kentcdodds

⁴¹ post by @kentcdodds

⁴² post by @dtcb

⁴³ post by @swyx

⁴⁴ OpenClaw’s Memory Sucks and the fix is simple — Dhravya Shah, Supermemory

⁴⁵ OpenClaw’s Memory Sucks and the fix is simple — Dhravya Shah, Supermemory

⁴⁶ OpenClaw’s Memory Sucks and the fix is simple — Dhravya Shah, Supermemory

⁴⁷ post by @steipete

⁴⁸ post by @steipete

⁴⁹ post by @steipete

⁵⁰ post by @karpathy

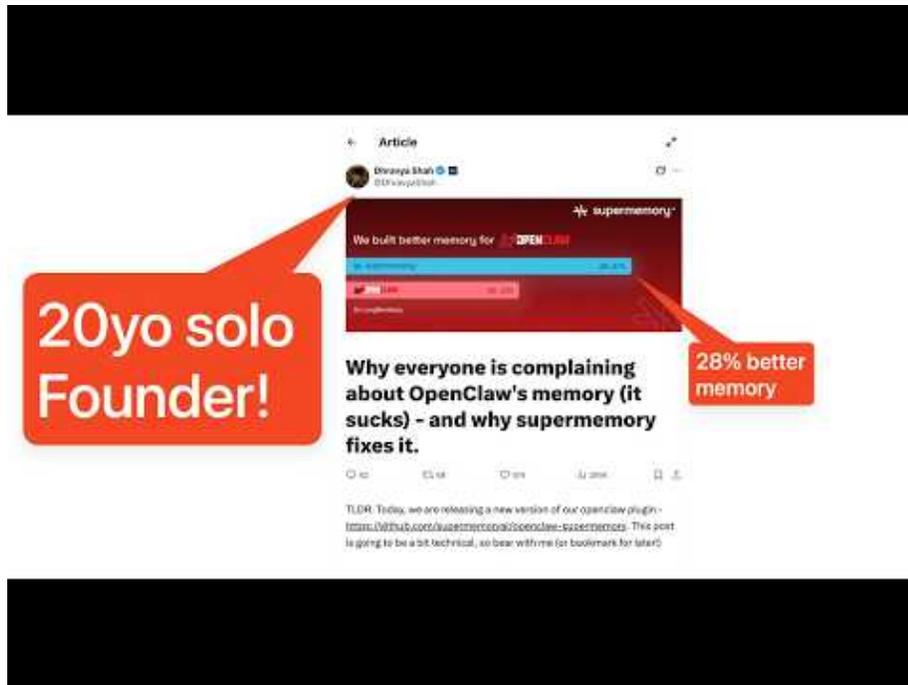
⁵¹ post by @karpathy

⁵² post by @theo

⁵³ post by @theo

WATCH & LISTEN

- **Latent Space** — **19:24–20:35**: why user profiles beat literal retrieval. Good explanation of why an agent needs a tiny always-on profile plus recent episodes to answer questions like what monitor fits you, even if you never explicitly talked about monitors ⁵⁴

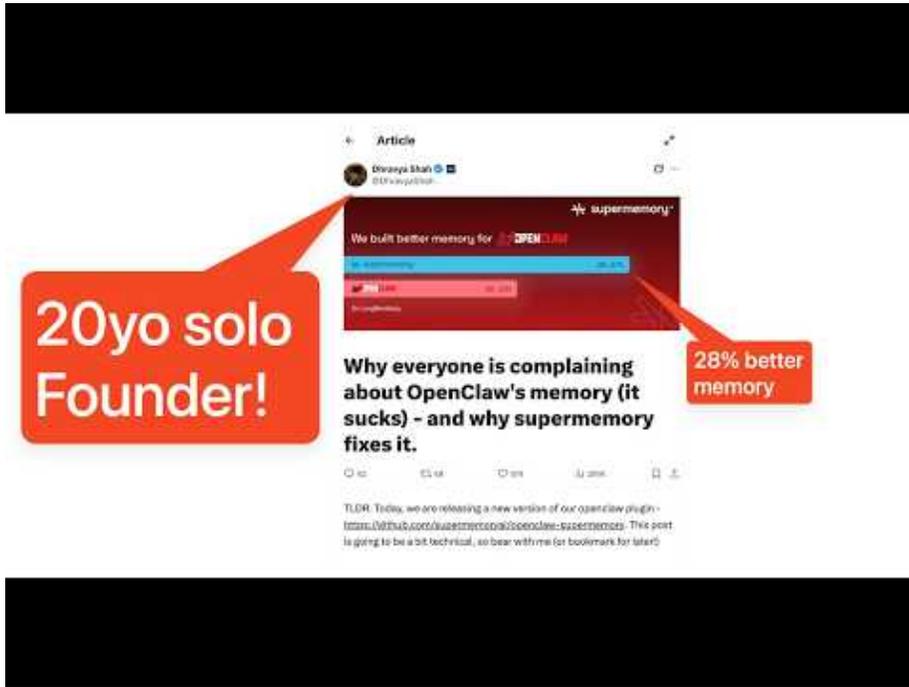


OpenClaw's Memory Sucks and the fix is simple — Dhruvya Shah, Supermemory (19:24)

- **Latent Space** — **22:25–23:42**: hybrid memory mode for OpenClaw. Memories surface first, RAG fills the gap when memory misses, and the system extracts that information in the background for future turns ⁵⁵

⁵⁴OpenClaw's Memory Sucks and the fix is simple — Dhruvya Shah, Supermemory

⁵⁵OpenClaw's Memory Sucks and the fix is simple — Dhruvya Shah, Supermemory



OpenClaw's Memory Sucks and the fix is simple — Dhruva Shah, Supermemory (22:25)

- **NVIDIA on Latent Space — 1:08:21–1:09:41:** why coding agents keep beating general agents. The argument is straightforward: the terminal gives agents access to compilers, tests, and every installed tool, so the feedback loop is tighter than pure chat ⁵⁶

⁵⁶NVIDIA's AI Engineers: Brev, Dynamo and Agent Inference at Planetary Scale and "Speed of Light"



NVIDIA's AI Engineers: Brev, Dynamo and Agent Inference at Planetary Scale and "Speed of Light" (68:21)

PROJECTS & REPOS

- **acpx** — bridge layer that lets Codex call OpenClaw over ACP and OpenClaw call Codex back. Steinberger is already using it for private agent-to-agent discussion with an approval gate before posting to Discord ⁵⁷⁵⁸⁵⁹⁶⁰
- **Super Memory** — open-source context infrastructure for stateful agents. Shah says the project reached **100k users** on about **\$5/month** of Cloudflare spend in its early consumer phase and hit **10k GitHub stars in a few weeks** after open source ⁶¹⁶²
- **Memory Benchmark** — open-source eval harness for memory systems across providers, benchmarks, and judges, with metrics for quality, latency, cost, top-K recall, and NDCG ⁶³
- **Karpathy's lightweight GitHub coordination pattern** — use Discussions for agent-written run summaries and PRs for exact commits you

⁵⁷ post by @steipete

⁵⁸ post by @steipete

⁵⁹ post by @steipete

⁶⁰ post by @steipete

⁶¹OpenClaw's Memory Sucks and the fix is simple — Dhruvya Shah, Supermemory

⁶²OpenClaw's Memory Sucks and the fix is simple — Dhruvya Shah, Supermemory

⁶³OpenClaw's Memory Sucks and the fix is simple — Dhruvya Shah, Supermemory

might adopt without merging ⁶⁴⁶⁵⁶⁶

Editorial take: the edge is shifting from choosing one best model to building the system around it — routing, memory, terminal access, and permission boundaries ⁶⁷⁶⁸⁶⁹⁷⁰

Sources

1. post by @dtcb
2. post by @swyx
3. NVIDIA's AI Engineers: Brev, Dynamo and Agent Inference at Planetary Scale and "Speed of Light"
4. post by @jediahkatz
5. post by @LangChain
6. post by @LangChain
7. OpenClaw's Memory Sucks and the fix is simple — Dhruvya Shah, Super-memory
8. post by @romainhuet
9. post by @steipete
10. post by @steipete
11. post by @steipete
12. post by @rileybrown
13. post by @kentcdodds
14. post by @kentcdodds
15. post by @karpathy
16. post by @theo

⁶⁴ post by @karpathy

⁶⁵ post by @karpathy

⁶⁶ post by @karpathy

⁶⁷ post by @swyx

⁶⁸ OpenClaw's Memory Sucks and the fix is simple — Dhruvya Shah, Supermemory

⁶⁹ NVIDIA's AI Engineers: Brev, Dynamo and Agent Inference at Planetary Scale and "Speed of Light"

⁷⁰ NVIDIA's AI Engineers: Brev, Dynamo and Agent Inference at Planetary Scale and "Speed of Light"