

Harness Engineering, Codex Playbooks, and Claude Code's Enterprise Signal

Coding Agents Alpha Tracker

2026-04-28

Harness Engineering, Codex Playbooks, and Claude Code's Enterprise Signal

By Coding Agents Alpha Tracker • April 28, 2026

DeepAgent's benchmark jump makes the case for harness tuning over model swapping. Also in today's brief: Riley Brown's Codex playbook, Applied Intuition's Claude Code adoption, GPT-5.5 workflow notes, and concrete patterns for context, evals, and automation.

TOP SIGNAL

Today's clearest edge is **harness engineering**, not another round of model swapping. Harrison Chase says LangChain moved **DeepAgent** from **30th to 5th** on Terminal Bench without changing the model, while Geoffrey Huntley's harness panel makes the same point from the builder side: exhaust tool design, context funneling, and evals before stacking more agent loops [1, 2].

Across the sources, the recurring pattern is that tool design, context management, and observability are producing gains people often chase via model changes alone [1].

TOOLS & MODELS

- **DeepAgent** — LangChain's open-source, model-agnostic harness jumped from **30th to 5th** on Terminal Bench via harness tuning alone. Good reminder: benchmark movement can come from the scaffold, not the weights [1].
- **Codex + GPT-5.5** — Riley Brown says his team of **seven engineers** has switched to Codex, and he prefers it to Claude Desktop and Cursor for complex infrastructure work because it combines coding, docs, automations, and multitasking chats in one GUI [3, 4].

- **GPT-5.5** — Greg Brockman says it is strong on hard tasks like writing GPU kernels. Riley’s usage note: use lower effort for simpler changes and higher effort for harder work, with the tradeoff that API usage costs more even if the model can be more direct about intent [5, 3].
- **GPT Image 2 -> Codex frontend loop** — Romain Huet points to a practical pairing: do the design pass with GPT Image 2, then build with Codex + GPT-5.5. The cited demo showed watch components generation and alignment landing in one shot [6, 7].
- **Claude Code at Applied Intuition** — Peter Ludwig says Cursor was initially the hottest tool internally, but Claude Code now leads their internal adoption leaderboard and is phenomenally useful. His caveat is the one that matters: in safety-critical systems, human validation is still mandatory [8].
- **OpenClaw 2026.4.26** — New release adds solutions for PR and issue management, remote test execution, and large CI testing workflows. The maintainer also calls out community work across docs, Windows, Linux, Docker, and local-model edge cases [9, 10].

WORKFLOWS & TRICKS

- **Codex multitasking setup**
 1. Create one project or folder per goal.
 2. Spawn separate chats with `Cmd+N`.
 3. Use the blue dot as the done or unread signal.
 4. Fork a chat when a branch deserves its own deliverable.
 5. If you want a Claude-specific pass, open the terminal with `Cmd+J` and run `claude` inside Codex [3, 11].
- **Reusable skill + examples beats long instructions** — Riley’s example is a YouTube research skill that pulls the last 10 transcripts and produces a targeted critique. He pairs that with a knowledge base of good examples, and says organized docs plus screen recordings are increasingly useful references for future agent workflows [3].
- **Automation recipe** — Do the task once, then tell the agent to turn it into a recurring automation. After that, use `run now`, inspect the output, and edit the automation until it works reliably [3].
- **Harness upgrade checklist** — Expose tools and knowledge through a virtual file-system-like interface when possible, use long-context models to reduce compaction, keep one context window per goal or activity, and only add more nested loops after you have exhausted tool design, system prompt work, context funneling, and evals [1, 2].
- **Trace-to-PR improvement loop** — Instrument traces, attach explicit or inferred failure signals, detect repeated tool loops, and let an agent suggest prompt, code, or harness changes from that data. Chase describes this as a meta-harness improvement flywheel, with a human deciding where the approval boundary sits [1].

PEOPLE TO WATCH

- **Harrison Chase** — High signal if you care about production agents, not demos. The useful bits today were harness engineering, trace-driven evals, and the idea that observability and regression testing remain timeless even as models improve [1].
- **Riley Brown** — Worth watching because he is publishing repeatable Codex workflows from real operating use, not just toy prompts. His recent output spans a full beginner guide and a Codex masterclass covering skills, browser agents, Claude Code inside Codex, and day-one projects [11, 12, 3].
- **Peter Ludwig** — Rare signal from a co-founder and CTO running a **1000-engineer** company in a safety-critical domain. His notes on tool adoption, hiring for AI engineer skills, and required human validation are more valuable than most benchmark chatter [8, 13].
- **Geoffrey Huntley’s harness panel** — Good source for timeless primitives: harness as deterministic code around the agent loop, nested loops as orchestration, and the warning not to stack more loops before fixing tool design and context handling [2].
- **Simon Willison** — Still one of the better voices for zooming out. His current framing is that agent loops may apply to more knowledge work than coding, and he says about **95%** of the code he produces today is not typed by him [14].

WATCH & LISTEN

- **16:30-18:31** — **Trace data into agent improvement.** Chase walks through the loop of running agents on data, evaluating failures, feeding the results back into the system, and deciding where the human should approve changes [1].



Automating the SDLC with LangChain, LangSmith, and Gemini (16:30)

- **33:21-36:59** — **Why examples beat instructions.** Riley explains why collecting strong outputs, organizing them in a knowledge base, and recording workflows can make agents much more reliable on subjective work [3].



Stop using Claude. Start using Codex? (33:21)

- **29:46-31:30** — **The simplest automation pattern in Codex.** Riley's recipe is dead simple: do the task once, convert it into an automation, test with `run now`, then edit until it is solid [3].



Stop using Claude. Start using Codex? (29:46)

PROJECTS & REPOS

- **DeepAgent** — The best project-level signal in the notes. Open-source, model-agnostic, and already showing that harness tuning alone can move a coding benchmark materially [1].
- **LangGraph + LangSmith** — Worth studying as a production stack for deterministic workflows, stateful resumes after failures, parallel execution, traces, online evals, and regression testing [1].
- **OpenClaw 2026.4.26** — Fresh release with PR and issue management, remote test execution, and heavier CI support, plus visible community contribution across platform edge cases [9, 10].

Editorial take: the durable edge right now is not one magic model pick; it is the harness and workflow around the model — cleaner contexts, reusable skills, trace-based evals, and human checkpoints where failure is expensive [1, 2, 8]

Sources

1. Automating the SDLC with LangChain, LangSmith, and Gemini
2. Context Engineering vs Harness Engineering vs Software Engineering
3. Stop using Claude. Start using Codex?

4. Meet The AI Creator Who Doesn't Write With AI
5. X post by @gdb
6. X post by @romainhuet
7. X post by @kagigz
8. The \$15B Physical AI Company: Simulation, Autonomy OS, Neural Sim, & 1K Engineers—Applied Intuition
9. X post by @openclaw
10. X post by @steipete
11. X post by @rileybrown
12. X post by @gregisenberg
13. Physical AI that Moves the World — Qasar Younis & Peter Ludwig, Applied Intuition
14. microsoft/VibeVoice