

OpenAI Frontier’s Harness Playbook, Typed Execution Layers, and Async Subagents

Coding Agents Alpha Tracker

2026-04-08

OpenAI Frontier’s Harness Playbook, Typed Execution Layers, and Async Subagents

By Coding Agents Alpha Tracker • April 8, 2026

OpenAI Frontier’s 1M+ LOC Codex experiment is today’s clearest practical signal: better coding-agent results are coming from harness design, not just model upgrades. Also: Deep Agents v0.5, OpenClaw’s latest release, GLM-5.1’s first real tests, and concrete workflows for specs, context hygiene, and feedback distillation.

TOP SIGNAL

OpenAI Frontier just published the clearest production harness playbook yet for coding agents: Ryan Lopopolo says his team built an internal beta over five months on a 1M LOC codebase with zero human-written code and no human-reviewed code pre-merge, using Codex across thousands of PRs [1]. The practical pattern is a harness around the model: sub-minute builds, agent-booted local stacks, markdown specs/skills/trackers, and agents that handle PR landing, CI flakes, and merge conflicts end-to-end [1, 2]. In Ryan’s framing, the scarce resource is now synchronous human attention—not tokens [1].

TOOLS & MODELS

- **Symphony** — OpenAI open-sourced Symphony, a “ghost library” and reference Elixir implementation for multi-agent Codex orchestration [1]. Lopopolo says Elixir was chosen because BEAM process supervision and GenServers map cleanly to task orchestration; the system was built after PR throughput jumped from ~3.5 to 5-10 PRs per engineer per day and humans became the context-switch bottleneck [2].
- **Deep Agents v0.5** — LangChain added remote **async subagents** that return a task ID immediately, stay stateful across follow-ups, and expose

`start/check/update/cancel/list` task ops [3]. They chose **Agent Protocol** over ACP/A2A for this release because it fits thread/run semantics and remote async work better right now [3].

- **OpenClaw + provider churn** — Anthropic no longer allows third-party harnesses like OpenClaw to use Claude cloud subscriptions the old way, forcing either extra usage or another provider [4]. Matthew Berman says swapping OpenClaw to **GPT-5.4** took ~3 minutes, but he also recommends model-specific prompt files because prompts that work for Opus and GPT diverge materially [4]. Meanwhile, **OpenClaw v2026.4.7** shipped `infer`, session branch/restore, webhook TaskFlows, and `memory-wiki` [5].
- **GLM-5.1** — Z.ai’s new open-weight model is **754B params / 1.51TB**, positioned for long-horizon tasks, with claims of #1 open-source and #3 global performance across SWE-Bench Pro, Terminal-Bench, and NL2Repo [6, 7]. More useful than the benchmark slide: Simon Willison ran it through `llm/OpenRouter` on a real SVG+animation task, then got a clean diagnosis and fix for a broken CSS/SVG transform bug on the next turn [8].
- **Claude Code /powerup** — Small but practical: Claude Code now has an interactive onboarding flow with **10 short lessons/demos**. Update with `claude update`, then run `/powerup` [9, 10].

WORKFLOWS & TRICKS

- **Steal Frontier’s automation loop:** 1) when the agent fails, add missing capability/context/structure instead of just re-prompting; 2) force builds under one minute; 3) let the agent boot the local stack via skills/scripts; 4) give it a landing skill that waits for reviewers/CI, fixes flakes/conflicts, and merges; 5) keep humans on release branches and smoke tests [1, 2].
- **Run spec-first, fresh-thread workflows:** photo the whiteboard or upload wireframes/PDFs, ask the agent clarifying questions to turn them into a spec, then open a new thread for each independent feature [11]. Huet’s warning is concrete: giant conversations plus the wrong folder scope create context overload, blank screens, and bad commits; use **GPT-5.4** for heavier scaffolding and **Codex Spark** for fast UI passes [11].
- **Stop pasting repos into context:** Ryan says models “crave text,” so Frontier keeps turning more state into text and local tooling the agent can inspect [2]. Theo’s version of the same rule: fetch the 8 relevant lines with `grep` instead of dumping 100k tokens into chat, because smaller deterministic retrieval beats bloated context on both cost and output quality [12].

“Agents are good at bash. Bash is not good for agents.” [13]

- **Treat feedback as data, not conversation:** Frontier stores session logs, PR comments, failed builds, and even Slack fixes, then feeds them back into skills/docs/tests so the whole team benefits [2]. LangChain frames

the same loop more generally: trace what the agent actually did, fix the harness, then measure whether the fix worked [14]. When writing skills, keep them narrow, make the description explicit about when to trigger, explain *why*, include examples, and push bulky refs/scripts into separate files [15].

PEOPLE TO WATCH

- **Ryan Lopopolo** — best production download of the day. The useful specifics are sub-minute builds, ghost libraries, inline dependencies, PR automation, and markdown self-review trackers—not just the 1B-token headline [1].
- **Theo Browne** — high-signal contrarian on execution layers. His argument: bash unlocked agents, but typed JS/TS runtimes are better long-term for permissions, isolation, approvals, and lower-context tool use; in the Cloudflare-style example he cites, token use dropped from 43.5k to 27k with better latency and a benchmark bump from 25.6 to 28.5 [13, 12].
- **Romain Huet** — worth tracking because he keeps the advice operational: fire parallel tasks, use worktrees, start fresh threads, generate specs from whiteboards, and reach for Spark when UI iteration speed matters more than deep reasoning [11].
- **Simon Willison** — still the fastest reality check on both new models and security impact. He showed GLM-5.1 fixing a real SVG/CSS bug in a follow-up prompt, and separately highlighted Anthropic’s Project Glasswing plus Nicholas Carlini’s claim that frontier models are now finding real, patched vulnerabilities at scale—including a 27-year-old OpenBSD issue [8, 16].
- **Peter Steinberger** — still worth following if you care about the control plane around agents: same-day OpenClaw and CodexBar ships, persistent memory, TaskFlows, and spend visibility across 16 providers [17, 5, 18].

WATCH & LISTEN

- **10:59-11:58** — “**Spawn the agent first.**” Ryan explains the inversion: don’t pre-bake the environment and then drop the agent in; start with Codex and give it the skills/scripts to boot the stack itself [2].



Extreme Harness Engineering for the 1B token/day Dark Factory — Ryan Lopopolo, OpenAI Frontier (10:59)

- **12:39-13:42** — **Markdown tables as control surfaces.** Tiny scaffolds like a tech-debt tracker and quality score let the agent audit business logic against written guardrails and propose its own follow-up work [2].



Extreme Harness Engineering for the 1B token/day Dark Factory — Ryan Lopopolo, OpenAI Frontier (12:39)

- **2:06-3:18 — Codex Spark for UI speed.** Huet shows the fast loop: generate a simple game, pin the preview, then iterate visual changes in seconds. Good calibration on when Spark beats a heavier model loop [11].



An app to bring people together · Web Dev Challenge S3.E3 (2:06)

- **21:53-22:43** — **Context overload in the wild.** Huet diagnoses a broken session as too much thread history plus the wrong folder scope—exactly the kind of setup bug that makes agents look worse than they are [11].



An app to bring people together · Web Dev Challenge S3.E3 (21:53)

PROJECTS & REPOS

- **Symphony** — OpenAI’s open-source orchestration blueprint for multi-agent Codex workflows. Adoption signal: it came out of a Frontier team working on a 1M LOC codebase, thousands of PRs, and reporting 5-10 PRs per engineer per day at peak [1, 2].
- **Agent Protocol + Deep Agents v0.5** — Remote, stateful async sub-agents with Python/JS example servers. Useful if you want background workers instead of inline blocking subagents [3].
- **OpenClaw v2026.4.7** — Added `infer`, webhook TaskFlows, session branch/restore, and `memory-wiki`; the release leans hard into persistent knowledge and multi-model workflows [5].
- **CodexBar 0.20** — Meta-tool for the multi-provider era: 16 providers tracked, new Perplexity/OpenCode Go backends, better account switching, and cleaner cost history [18].
- **GLM-5.1 weights** — 754B open weights positioned for long-horizon coding tasks; Simon’s real-world test was not a benchmark screenshot but a broken SVG animation that the model then diagnosed and fixed on the next turn [8, 7].

Editorial take: today’s edge is coming from harness engineering, not just model shopping—tight context, fast build/feedback loops, traceable failures, and exe-

cution layers built for what agents can actually read and control [1, 14, 12, 2].

Sources

1. Extreme Harness Engineering for Token Billionaires: 1M LOC, 1B toks/day, 0% human code, 0% human review — Ryan Lopopolo, OpenAI Frontier & Symphony
2. Extreme Harness Engineering for the 1B token/day Dark Factory — Ryan Lopopolo, OpenAI Frontier
3. Deep Agents v0.5
4. Anthropic banned OpenClaw...
5. X post by @openclaw
6. X post by @simonw
7. X post by @Zai_org
8. GLM-5.1: Towards Long-Horizon Tasks
9. X post by @morganlunt
10. X post by @_catwu
11. An app to bring people together · Web Dev Challenge S3.E3
12. The language holding our agents back.
13. X post by @theo
14. X post by @LangChain
15. The Complete Guide to Creating and Using Claude Skills 2026
16. Anthropic's Project Glasswing - restricting Claude Mythos to security researchers - sounds necessary to me
17. X post by @steipete
18. X post by @steipete