

# Verification becomes the bottleneck as Codex/Cursor/Claude push longer-running coding agents

Coding Agents Alpha Tracker

2026-02-19

## Verification becomes the bottleneck as Codex/Cursor/Claude push longer-running coding agents

*By Coding Agents Alpha Tracker • February 19, 2026*

Codex app workflows point to the new bottleneck: verification—teams are moving from “review code” to “review evidence.” Plus: SWE-bench’s independently-run Feb 2026 leaderboard, OneContext’s git-like cross-agent memory, Cursor’s past-conversation context, and new data on real-world Claude Code autonomy/oversight.

### TOP SIGNAL

OpenAI’s **Codex app** is pushing a clear practitioner lesson: as agents generate *more* code, **verification becomes the bottleneck**—not implementation speed. In a team discussion, Codex builders describe getting to “too much code to review,” and shifting toward evidence-driven PRs (agents clicking through flows, screenshotting, and uploading proof) as a way to reduce “code as proxy” verification [1].

### TOOLS & MODELS

- **OpenAI Codex app + GPT-5.3-Codex**
  - The Codex app launched last week and hit **1M+ downloads in the first week** [1].
  - OpenAI’s Thibault Sottiaux says **GPT-5.3-Codex** “beats every single other model” on top coding benchmarks [1] and that the team is still experimenting with how to supervise increasingly capable (and eventually multi-agent) systems [1].

- Product direction: a dedicated “**command center**” UI optimized for steering/supervision (vs. doing tasks yourself) [1], with features like **voice prompting, diagrams, and image rendering** [1].
- Latency work: they rewrote serving around **websockets + persistent connections**, targeting **~30–40% lower turn latency**, and plan additional optimizations to make the experience **2–3× faster** [1].
- **Codex is (still) open source + app-server is the integration surface**
  - Romain Huet highlights that the **Codex agent is open source**, including the **app server**, and it’s the same server behind the Codex app and integrations like Xcode [2, 3].
  - OpenAI devs also emphasize the shared open-source interface **app-server** for embedding Codex into products (including sign-in with ChatGPT) [4, 5]. Docs: <https://developers.openai.com/codex/app-server/> [3].
- **SWE-bench (Feb 2026) leaderboard update (independent run)**
  - SWE-bench published fresh “**Bash Only**” results run against the current generation of models using **mini-swe-agent** (~9k lines Python) and **one shared system prompt** across models for comparability [6].
  - Top scores (“% Resolved”) include **Claude 4.5 Opus 76.8%** and **Gemini 3 Flash 75.8%**; OpenAI’s **GPT-5.2** appears at **72.8%** [6].
  - Notably absent: **GPT-5.3-Codex**, which Simon Willison suggests may be because it isn’t available in the OpenAI API [6, 7].
- **Cursor 2.5: past conversations as context**
  - Cursor says it can now **use past conversations as context** for continuity across dev sessions [8]. Changelog: <http://cursor.com/changelog/2-5> [9].
- **OneContext (new): Git-like, cross-agent persistent memory**
  - OneContext is a new CLI that implements a “**Git Context Controller**” approach: memory stored as a simple file hierarchy (e.g., `main.md`, branch folders, commit logs) and reused across **sessions and coding agents** [10].
  - Install + run: `npm i -g onecontext-ai` then `onecontext` [10].
  - Motivation claim: despite 1M-token windows, “effective” context is closer to **120–200k**, so compaction/memory structures matter [10].
- **Anthropic: measuring agent autonomy in practice (usage telemetry)**
  - Anthropic analyzed millions of Claude Code + API interactions to characterize autonomy/oversight in the wild [11].
  - Key numbers: median Claude Code turn **~45s**, but the **99.9th percentile** duration grew from **<25 min to >45 min** over three months [12].
  - Oversight shift: after ~750 sessions, **>40%** of sessions are fully auto-approved; interruptions also rise with experience (5% → 9%) [13,

14].

- Risk surface: **73%** of tool calls appear human-in-the-loop; **0.8%** are irreversible; frontier cases include security systems, financial transactions, and production deployments [15].
- **LangSmith Agent Builder update + “evals as code” case study**
  - LangChain shipped updates to **LangSmith Agent Builder**: always-available agent chat, one-click “Chat → Agent,” file uploads, and a tool registry [16].
  - monday.com describes a code-first eval strategy for production service agents, claiming **8.7× faster** eval feedback loops (162s → 18s) by combining **parallelized Vitest** and **concurrent LLM evals**, plus GitOps-style CI/CD for evaluation logic [17].

## WORKFLOWS & TRICKS

- **Shift PR verification from “review code” to “review evidence” (Codex app pattern)**
  - The Codex team describes a pragmatic move: have the agent **run the app, click around, take screenshots for evidence, and upload proof to the PR**—so reviewers verify behavior directly rather than treating code review as the only proxy for correctness [1].
- **Use UI features to keep supervision tractable as agents parallelize**
  - Codex builders argue that as you run **many tasks in parallel**, a dedicated UI matters for “visibility into what the system is doing” and for steering/supervising outcomes [1].
  - Concrete interaction primitives:
    - \* **Mid-turn steering**: send new instructions while the model is still working; it adapts in real time [1].
    - \* **Plan mode**: quick Q&A plus an editable plan [1].
    - \* **Review mode**: annotate diffs with findings/stylistic notes [1].
- **Automations that actually ship: treat agents like scheduled workers**
  - Keep PRs mergeable: run a time-scheduled automation that **resolves merge conflicts** and **fixes build issues** via GitHub/CI skills (hourly / every two hours) [1].
  - “Random file” bug hunting: periodically pick a random file (via `Python rand`), find a subtle bug, fix and merge—reported to catch latent issues that wouldn’t be found otherwise [1].
  - Daily “what merged” report: get a themed summary of merged contributions to stay oriented during chaotic pre-launch periods [1].
- **Memory as a repo: branch/commit/merge your agent’s context (OneContext pattern)**
  - Minimal structure:
    - \* `main.md` for global roadmap/context [10]
    - \* per-branch folders with `commit.md` (milestones) and `log.md` (raw

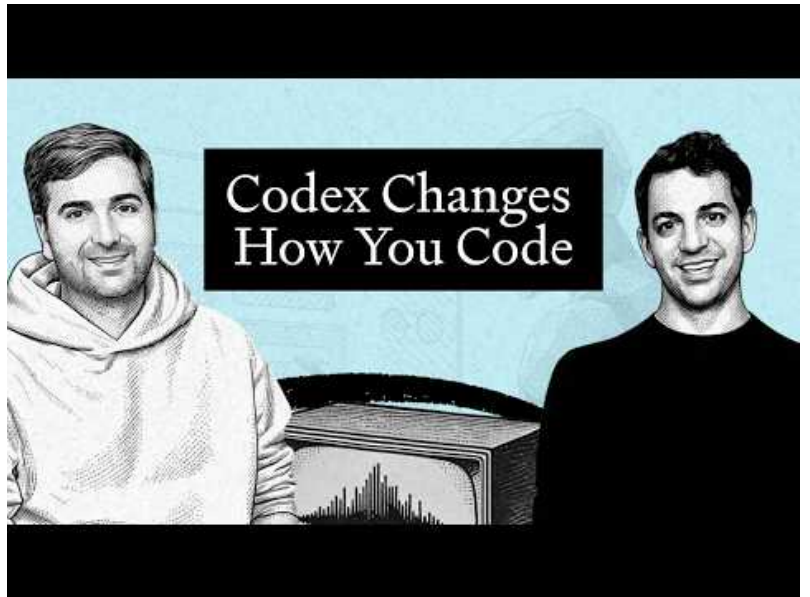
- conversations), plus metadata [10]
  - Three actions to make this operational for long tasks:
    1. **branch** when exploring an alternative strategy [10]
    2. **commit** at milestones/subtasks to summarize progress [10]
    3. **merge** when an approach is validated to roll learnings into main [10]
- **Pre-push agent review (anti-slop hygiene)**
  - Armin Ronacher’s guidance if you let an agent commit: run `/review` in a **fresh empty session locally** until the agent is satisfied *before* pushing; he calls PR machine-based review “horrible” as a default workflow [18].

## PEOPLE TO WATCH

- **Thibault Sottiaux (OpenAI Codex)** — high-signal details on super-*vision*/UI, speed/latency engineering, and where the bottleneck is moving (verification) [1].
- **Romain Huet + Alexander Embiricos (OpenAI)** — consistently pointing builders to the practical integration surface: the open-source Codex agent + **app-server** you can embed in products [2, 4].
- **Simon Willison** — good benchmark hygiene + practical browser automation: used **Claude for Chrome** to inject JS into SWE-bench charts to render percentage labels for clearer screenshots [6].
- **Jason Zhou** — clear demo of file-structured, cross-session memory and multi-agent “shared context” workflows using OneContext [10].
- **Addy Osmani** — pragmatic enterprise framing: the hard problem is **coordination**, not generation; focus on orchestrating a modest set of agents with control/traceability [19].

## WATCH & LISTEN

- **Codex app: “verification is the bottleneck now” + evidence-backed PRs** (~36:41–41:19)
  - Hook: why reviewers are drowning in generated code, and what “reviewing evidence” looks like (agents click through flows and attach



screenshots to PRs). [1]

*OpenAI's Codex: This Model Is So Fast It Changes How You Code (36:41)*

- **OneContext demo: cross-folder memory + stop-hook summaries** (~10:35–14:21)
  - Hook: a concrete walkthrough of persistent memory across sessions/agents, with post-session logging + summarization via a “stop



hook.” [10]

*Agent memory resolved? (10:35)*

## PROJECTS & REPOS

- **openai/codex (open source)** — repo link: <https://github.com/openai/codex> [3].
- **Codex app-server docs** — integration entrypoint for embedding Codex into your app: <https://developers.openai.com/codex/app-server/> [3].
- **OneContext CLI** — installable tool packaging the Git-style context controller approach (`npm i -g onecontext-ai`) [10].
- **mini-swe-agent prompts** — SWE-bench harness prompts are public: [https://github.com/SWE-agent/mini-swe-agent/blob/v2.2.1/src/minisweagent/config/benchmarks/swebench\\_prompts.txt](https://github.com/SWE-agent/mini-swe-agent/blob/v2.2.1/src/minisweagent/config/benchmarks/swebench_prompts.txt) [6].

— **Editorial take:** Agents are getting fast enough that *shipping* is now gated by **verification, memory hygiene, and eval discipline**—not raw code generation [1, 10, 17].

---

### Sources

1. OpenAI's Codex: This Model Is So Fast It Changes How You Code
2. X post by @romainhuet
3. X post by @romainhuet
4. X post by @embirico
5. X post by @reach\_vb
6. SWE-bench February 2026 leaderboard update
7. X post by @simonw
8. X post by @cursor\_ai
9. X post by @cursor\_ai
10. Agent memory resolved?
11. X post by @AnthropicAI
12. X post by @AnthropicAI
13. X post by @AnthropicAI
14. X post by @AnthropicAI
15. X post by @AnthropicAI
16. X post by @LangChain
17. monday Service + LangSmith: Building a Code-First Evaluation Strategy from Day 1
18. X post by @mitsuhiko
19. How to Build Reliable AI at Scale: Insights from Addy Osmani