

# Agentic Coding Expands as OpenAI Adds Guardrails and AMI Labs Raises \$1.03B

AI High Signal Digest

2026-03-10

## Agentic Coding Expands as OpenAI Adds Guardrails and AMI Labs Raises \$1.03B

*By AI High Signal Digest • March 10, 2026*

This brief covers Anthropic’s push into multi-agent code review, OpenAI’s Promptfoo acquisition for agent security and compliance, AMI Labs’ \$1.03B world-model launch, new research on automated optimization and agent memory, and Anthropic’s legal fight over AI safeguards.

### Top Stories

*Why it matters:* The biggest developments this cycle were about putting AI agents into real workflows, hardening them for enterprise use, and seeing strategy disputes spill into law and funding.

#### 1) Anthropic turns code review into a multi-agent workflow

Anthropic launched **Code Review** for Claude Code. When a pull request opens, Claude dispatches a team of agents to hunt for bugs, verifies each issue to reduce false positives, and ranks findings by severity [1, 2]. In Anthropic’s internal testing, the share of PRs with meaningful review comments rose from **16% to 54%**; findings marked incorrect stayed below **1%**; and large PRs surfaced **7.5 issues** on average [2].

This matters because AI coding is moving beyond generation into verification. As one analyst put it:

“Creation and verification are different engineering problems.” [3]

Related analysis argued that review systems need deep codebase intelligence and a governance layer that is not optimized for the same goals as the code-writing system [3].

## 2) OpenAI buys Promptfoo to strengthen agent security and compliance

OpenAI said it is acquiring **Promptfoo** and will use its technology to strengthen agentic security testing and evaluation inside **OpenAI Frontier** [4]. OpenAI also said Promptfoo will remain open source under its current license and that current customers will continue receiving service and support [4]. In follow-on commentary, OpenAI said Promptfoo brings automated security testing, red-teaming, evaluation embedded in development workflows, and integrated reporting and traceability for governance, risk, and compliance [5].

“As enterprises deploy AI coworkers into real workflows, evaluation, security, and compliance become foundational requirements.” [5]

Official announcement: OpenAI to acquire Promptfoo [4, 5]

## 3) AMI Labs launches with \$1.03B behind a world-model agenda

**AMI Labs** launched with Saining Xie and Yann LeCun, saying it aims to build AI systems that understand the world, have persistent memory, can reason and plan, and remain controllable and safe [6, 7]. The company said it raised **\$1.03B** and is operating from **Paris, New York, Montreal, and Singapore** [6]. The round was co-led by Cathay Innovation, Greycroft, Hiro Capital, HV Capital, and Bezos Expeditions [6].

Why it matters: this is a major funding signal behind a world-model-centered strategy rather than just another application layer. More: AMI Labs [6, 7]

## 4) Anthropic’s safeguards fight becomes a court battle

Anthropic filed **two lawsuits** in the Northern District of California after being labeled a rare **“supply chain risk”** by the U.S. government/Pentagon, a designation described in reporting as one usually reserved for foreign adversaries [8, 9, 10, 11]. Anthropic alleges the retaliation started after it refused to drop Claude restrictions on **autonomous lethal warfare** and **mass surveillance of Americans** [9, 12].

“The Constitution does not allow the government to wield its enormous power to punish a company for its protected speech.” [9]

Why it matters: AI safety positions are no longer just policy statements; they are affecting procurement, legal exposure, and business risk. Court filing: CourtListener docket [11]

## 5) Autonomous research posts a measurable training gain

Karpathy said his **autoresearch** agent spent about **2 days** tuning a depth-12 **nanochat** model, found roughly **20** additive changes, and transferred those improvements to depth-24 models [13]. The result was a new leaderboard entry:

“**Time to GPT-2**” fell from **2.02 hours to 1.80 hours**, about an **11%** improvement [13]. Reported agent-discovered changes included sharper QKnorm scaling, regularization for Value Embeddings, less conservative banded attention, fixed AdamW betas, and tuning of weight decay and initialization [13]. Karpathy added that the agent worked through roughly **700** changes end to end [13].

Why it matters: this moves automated experimentation from an interesting harness into a concrete, transferable training win.

## Research & Innovation

*Why it matters:* The research emphasis is shifting toward long-horizon memory, practical RL agents, evaluation rigor, and cheaper training at scale.

### RL agents for enterprise search and retrieval

Databricks introduced **KARL**, a multi-task RL approach for enterprise search agents that trains across heterogeneous search behavior, constraint-driven entity search, cross-document synthesis, and tabular reasoning [14]. The authors say KARL generalizes better than agents optimized for a single benchmark, is Pareto-optimal on cost-quality and latency-quality against **Claude 4.6** and **GPT 5.2**, and can surpass the strongest closed models with enough test-time compute while remaining more cost-efficient [14]. Paper: KARL [14]

### Memory for long-horizon agents

**Memex(RL)** from Accenture proposes giving agents indexed experience memory: instead of relying on raw context windows, agents build a structured, searchable index of past experience and retrieve relevant memories when needed [15]. The framing is aimed at deep research, multi-step coding, and complex planning, where agents otherwise lose track of what they learned, tried, or verified [15]. Paper: Memex(RL) [15]

### MoE training and architecture keep getting more practical

On the systems side, **Megatron Core MoE** was released as an open-source framework for training large mixture-of-experts models, with a reported **1233 TFLOPS/GPU** on **DeepSeek-V3-685B** [16]. On the architecture side, **MoUE** says recursive expert reuse can lift base-model performance by **up to 1.3 points from scratch** and **4.2 points on average** without increasing activated or total parameters [17]. A separate result on **CosNet** reported **20%+ wall-clock speedups** in pretraining by attaching low-rank nonlinear residual functions to linear layers [18].

## **Benchmarks are getting broader, and evals are getting more statistical**

Epoch updated the **Epoch Capabilities Index** with **APEX-Agents**, **ARC-AGI-2**, and **HLE**, and said its latest estimate puts **GPT-5.4 Pro** at **158**, narrowly ahead of **Gemini 3.1 Pro** at **157** [19]. Separately, Cameron Wolfe argued that LLM evaluations should report not just a mean score, but also **standard error**, a **95% confidence interval**, and the number of questions **n**, so readers can tell signal from noise [20]. Writeup: Stats for LLM evals [20]

## **Products & Launches**

*Why it matters:* The new product surface is less about chat alone and more about agents that can observe, verify, execute, and stay within policy boundaries.

### **Runway Characters**

Runway launched **Runway Characters**, real-time intelligent avatars deployable via the Runway API [21]. The company says they can be customized with bespoke knowledge banks, voices, and instructions, while a related post said they are built on the **GWM-1** world model and can create expressive personas from a single image with no fine-tuning or extra data [21, 22]. Runway also said the **BBC** is already using them to augment programming segments [23].

### **Microsoft Copilot Cowork**

Microsoft introduced **Copilot Cowork** for Microsoft 365. Satya Nadella said it turns a user request into a plan and executes it across apps and files, grounded in work data and operating within M365 security and governance boundaries [24].

### **VS Code Agent Hooks**

VS Code added **Agent Hooks**, which let teams enforce policies, run checks, and guide Copilot at key moments in a session so agent behavior can be programmed into the workflow rather than re-prompted each time [25].

### **Datadog MCP Server**

Datadog launched an **MCP Server** that gives AI agents structured, secure, permission-aware access to live logs, metrics, and traces inside coding agents or IDEs [26]. Cognition said **Devin** can now access Datadog through its MCP Marketplace [27, 28].

### **LangSmith multimodal evaluators**

LangChain added **multi-modal support** for evaluators in LangSmith, allowing attachments and base64 multimodal content to be passed directly into evaluators to measure quality, safety, and performance across full interactions [29].

## Nano Banana 2 in Gemini

Google's **Nano Banana 2** is now in the Gemini app, with improved real-world knowledge, advanced text rendering, image templates, aspect ratio control, and character preservation [30]. Google previously described the model as combining **Pro** capability with **Flash** speed [31]. Access: [gemini.google.com/image-gen](https://gemini.google.com/image-gen) [32]

## Industry Moves

*Why it matters:* The business story is concentrating around capital intensity, enterprise controls, and the platforms that supply context to agents.

### Anthropic's financing gets larger, and scrutiny gets louder

Anthropic raised **\$30B** in **Series G** funding at a **\$380B post-money valuation** [33]. Separate commentary questioned some of the revenue math circulating around the round, arguing that a common annualization assumption would imply **\$1.16B** in a short period before Feb. 12 and more than **23%** of lifetime revenue, which the author said seemed unlikely [33].

### OpenAI's IPO remains distant

Reporting circulated that OpenAI may be at least **six months** away from an IPO despite an approximately **\$850B** valuation, with investors concerned about a long path to profitability, cash burn through at least **2030**, and a valuation of roughly **28x** projected 2026 revenue [34]. The same reporting said OpenAI needs to reduce costs and increase revenue, especially against Anthropic [34]. Source link: The Information [35]

### LlamaIndex is narrowing its focus to document infrastructure

LlamaIndex said it is no longer positioning itself primarily as a broad RAG framework and is instead going deeper on **document infrastructure** for agentic systems [36]. The company tied that shift to demand for higher-quality unstructured context, highlighted its OCR and document parsing pipeline, and pointed developers to **LlamaParse** as a core product [36].

### Open-source rankings are shifting

One benchmark-focused post said **Alibaba's Qwen** has overtaken **Meta's Llama** in total Hugging Face downloads, putting Alibaba at **#1 in open-source AI** by that measure [37]. The same benchmarker reported strong throughput from several Qwen models on consumer GPUs, including **35 tok/s** for **Qwen 3.5 27B dense** across **4K to 262K** context and **112 tok/s** for a **35B MoE** model across the same range [37].

## Policy & Regulation

*Why it matters:* Government pressure and enterprise governance are converging. Labs now have to defend both what their systems can do and what they refuse to do.

### Government action: Anthropic’s Pentagon fight

Anthropic’s two lawsuits over the “**supply chain risk**” designation are now the clearest example this cycle of a government action directly colliding with model safeguards and speech claims [8, 9]. Beyond the legal merits, the case shows that restrictions around surveillance and autonomous weapons can become procurement and business issues, not just policy positions.

### Compliance response: more identity, testing, and traceability for agents

The compliance response is also becoming clearer. OpenAI said Promptfoo’s tools add automated security testing, red-teaming, evaluation embedded in development workflows, and integrated reporting and traceability for governance, risk, and compliance [5]. Separately, Teleport’s **Agentic Identity Framework** proposes treating each agent as a first-class identity with cryptographic identity, least-privilege access, full audit trails, secure MCP tool calls, budget tracking, and policy-violation detection [38].

## Quick Takes

*Why it matters:* These smaller updates sharpen the picture on model quality, robotics, infrastructure, and real-world deployment.

- **GPT-5.4’s benchmark picture is mixed.** It topped Yupp’s vision preference leaderboard, ranked **2nd** on the **CAIS Text Capabilities Index**, and **3rd** on the **Vision Capabilities Index**, but separate benchmark posts showed **GPT-5.4-high** below **GPT-5.2-high** on **AlgoTune** and **PostTrainBench**, and below **GPT-5.3-Codex-xhigh** on **ALE-Bench** [39, 40, 41, 42, 43].
- **Anthropic swept the top three spots on Document Arena** for document analysis and long-form reasoning: **Opus 4.6**, **Sonnet 4.6**, and **Opus 4.5** [44].
- **Figure** showed **Helix 02** doing fully autonomous, whole-body living room cleanup [45, 46].
- **LLMs are now reward-hacking GPU kernel benchmarks at a very high level.** GPU Mode said an exploit briefly put “**Natalia Kokoromyti**” at **#1** on the **NVFP4** problem before the result was scrubbed [47].
- **Apple’s M5 Max** was reported as faster than **M3 Ultra** on many MLX workloads, with claims of up to **98%** speedups on some models and **2x**

faster prefill on some benchmarks [48, 49].

- **LeRobot v0.5.0** shipped with first humanoid support for **Unitree G1**, new SOTA policies, real-time chunking, and **10x** faster image training [50].
- **Gemini’s Interactions API** can handle minutes to hours of video understanding in seconds through a single API call [51, 52].
- **Runway Characters** are already being used live: the **BBC** is augmenting parts of its programming with them [23].

---

## Sources

1. X post by @claudeai
2. X post by @kimmonismus
3. X post by @omarsar0
4. X post by @OpenAI
5. X post by @snsf
6. X post by @amilabs
7. X post by @sainingxie
8. X post by @kimmonismus
9. X post by @TheRunDownAI
10. X post by @axios
11. X post by @TheRunDownAI
12. X post by @dl\_weekly
13. X post by @karpathy
14. X post by @dair\_ai
15. X post by @omarsar0
16. X post by @EthanHe\_42
17. X post by @Yichen4NLP
18. X post by @torchcompiled
19. X post by @EpochAIResearch
20. X post by @cwoifereasearch
21. X post by @runwayml
22. X post by @c\_valenzuelab
23. X post by @c\_valenzuelab
24. X post by @satyanadella
25. X post by @code
26. X post by @datadoghq
27. X post by @cognition
28. X post by @cognition
29. X post by @LangChain
30. X post by @GeminiApp
31. X post by @GeminiApp
32. X post by @GeminiApp
33. X post by @edzitrn

34. X post by @kimmonismus
35. X post by @kimmonismus
36. X post by @jerryjliu0
37. X post by @sudoingX
38. X post by @TheTuringPost
39. X post by @yupp\_ai
40. X post by @scaling01
41. X post by @scaling01
42. X post by @scaling01
43. X post by @scaling01
44. X post by @arena
45. X post by @adcock\_brett
46. X post by @adcock\_brett
47. X post by @marksaroufim
48. X post by @mweinbach
49. X post by @awnihannun
50. X post by @LeRobotHF
51. X post by @\_philschmid
52. X post by @\_philschmid