

ER Diagnostic AI, Paid Agent Workflows, and Open-Model Security Risks

AI High Signal Digest

2026-05-11

ER Diagnostic AI, Paid Agent Workflows, and Open-Model Security Risks

By AI High Signal Digest • May 11, 2026

A new ER diagnosis study put an older OpenAI model ahead of physicians in early triage, while Codex showed a fully autonomous path from prompt to payment. This brief also covers new alignment and tool-use research, an open-model data extraction risk, notable product releases, and fresh moves from Ramp, China Mobile, and MiniMax.

Top Stories

Why it matters: These were the clearest real-world signals on where AI is gaining capability, and where that capability could matter quickly.

- **A new Science study found OpenAI's o1 outperforming ER physicians on diagnosis.** The model reached **67%** correct or near-correct diagnoses versus **50–55%** for doctors, with the widest gap appearing in early triage when information is limited [1]. The same writeup said o1 was near-perfect on structured clinical reasoning, but the study covered only short ER encounters and did not test imaging [1].
- **Codex completed a full bounty workflow and got paid.** In one public example, a user prompted Codex to “make me \$5”; it found an open-source security bounty, opened a legitimate PR, followed up with the maintainer, handled the verification loop, protected payment details, and earned **\$16.88** after about **22 hours** [2]. The poster estimated a **\$506.40/month** run rate if repeated daily, and Sam Altman called the example “interesting” [2, 3].

Research & Innovation

Why it matters: The most useful research this cycle targeted alignment, tool reliability, and a still-open security problem in model training data.

- **Model Spec Midtraining** cut agentic misalignment from **54%** to **7%**, outperforming deliberative alignment baselines [4].
- **Apple’s reviewer-agent paper** moves evaluation into the execution loop: a reviewer inspects provisional tool calls before they run and feeds back corrections [5]. Reported gains were **+5.5%** on BFCL irrelevance detection, **+1.6%** on relevance, and **+7.1%** on **2-Bench** multi-turn, all without retraining the base agent [5]. The paper also introduced Helpfulness-Harmfulness metrics and argued the reviewer can be optimized as a separate production lever [5].
- **A Google DeepMind ablation highlighted a data-extraction risk in open-weight models.** It found that prompting with only the **chat template** can cause models to regurgitate their **SFT** and even **RL** training data, including verbatim RL QA samples [6, 7]. Separate testing claimed the **Magpie** method still extracted DeepSeek SFT data with a specific prompt, surfacing mostly math problems and a file labeled **Communism_alignment.csv** [8, 9].

Products & Launches

Why it matters: New releases kept pushing on retrieval quality, multimodal generation, and longer-running agent workflows.

- **Qdrant 1.17** adds what it calls the first **vector index-native relevance feedback** approach, aiming to push relevance into retrieval itself for smarter vector search [10].
- **HiDream-O1-Image** launched on fal.ai with a unified **pixel-level transformer** that processes raw pixels, text, and task cues in one token space; fal highlights stronger long-text layouts and better subject consistency across scenes [11].
- The **Codex macOS app** now supports **long-running threads with heartbeats, automations**, and integrations with **GitHub, Gmail**, and more; users also said recent updates made it much faster [12].

Industry Moves

Why it matters: Companies are still testing whether advantage comes from custom post-training, aggregation layers, or inference efficiency.

- **Ramp Labs and PrimeIntellect** built **Fast Ask**, a small RL-trained subagent for spreadsheet questions that scored **+4% over Opus** on exact-match accuracy at **Haiku latency** [13].
- **China Mobile** launched **MoMa**, a MaaS platform integrating **300+ models**. It claims centralized token procurement cuts costs by **30%+**

and resource use by **50%+**, with **billion-level daily token calls** and plans starting at **¥5.99** [14]. One analyst argued it looks like a state-owned OpenRouter equivalent with limited differentiation [15].

- **MiniMax and NVIDIA** said they are deepening collaboration on **inference optimization** for next-generation models, and MiniMax previewed a new **sparse solution** coming soon [16].

Quick Takes

Why it matters: These smaller updates still sharpen the picture on science, open models, and coding performance.

- **University of Warwick’s RAVEN AI** scanned data across **2.2M stars**, confirming **118** new exoplanets and identifying **2,000+** candidates, nearly **1,000** previously unspotted [17].
- The **GGUF ecosystem** on Hugging Face reached **176,000** public models; monthly additions rose from about **5.1K** in Oct–Feb to about **9.2K** in March–April [18].
- The **Continuous Latent Diffusion Language Model** paper was released, with experiments reported to scale up to **2,000 EFLOPs** [19, 20].
- Independent testers called **GPT-5.5 high** the strongest coding agent they had measured, while also warning that reduced thinking budgets can hurt high-complexity bug-finding; another developer said it was the first frontier model to solve his long-running refactor test [21, 22, 23].

Sources

1. X post by @kimmonismus
2. X post by @chatgpt21
3. X post by @sama
4. X post by @TheAITimeline
5. X post by @omarsar0
6. X post by @Ex0byt
7. X post by @teortaxesTex
8. X post by @sheryuo
9. X post by @teortaxesTex
10. X post by @qdrant_engine
11. X post by @fal
12. X post by @reach_vb
13. X post by @RampLabs
14. X post by @yyyole
15. X post by @teortaxesTex
16. X post by @RyanLeeMiniMax
17. X post by @TheRunDownAI
18. X post by @ClementDelangue

19. X post by @_akhaliq
20. X post by @TheAITimeline
21. X post by @hiarun02
22. X post by @MParakhin
23. X post by @jamesjyu