

# Claude extraction allegations, video-trained action agents, and faster real-time AI APIs

AI High Signal Digest

2026-02-24

## Claude extraction allegations, video-trained action agents, and faster real-time AI APIs

*By AI High Signal Digest • February 24, 2026*

Anthropic says it uncovered industrial-scale Claude extraction campaigns, while new agent capabilities arrive via video-trained computer-action models (FDM-1) and OpenAI’s lower-latency WebSockets + upgraded real-time voice. Also: a notable interpretability release (Steering-8B), DeepSeek V4 signals amid export-ban reporting, and a broad set of product and enterprise moves.

### Top Stories

#### 1) Anthropic reports industrial-scale “distillation” campaigns targeting Claude

*Why it matters:* If model capabilities can be reconstructed through massive API querying, it changes the security perimeter from *compute access* to *API abuse detection*—with implications for safeguards, export controls, and the economics of frontier development.

Anthropic says it identified industrial-scale campaigns by **DeepSeek, Moonshot AI, and MiniMax** that used **~24,000 fraudulent accounts** to generate **16M+ exchanges** with Claude, extracting capabilities to train or improve other models<sup>12</sup>. One breakdown circulating in the discussion lists query volume as **~150k (DeepSeek), 3.4M (Moonshot), and 13M (MiniMax)**<sup>3</sup>.

A widely shared explanation frames distillation as copying a stronger model’s behavior by collecting millions of input/output examples—common internally for “smaller, cheaper models,” but here characterized as an espionage-like operation

---

<sup>1</sup> post by @AnthropicAI

<sup>2</sup> post by @niubi

<sup>3</sup> post by @Bayesian0\_0

via API calls <sup>45</sup>. Specific tactics cited include requesting step-by-step reasoning, targeting agent capabilities across many accounts, and rapidly pivoting to new model releases <sup>6</sup>.

Several posts emphasized downstream risks: safety filters can be lost, and “frontier AI without safeguards” could be used in military/surveillance contexts <sup>789</sup>. Others argue this could undermine chip export controls if output-copying scales <sup>1011</sup>. A recurring recommendation: **shared detection systems across frontier labs**, since attackers can rotate to the weakest defenses <sup>12</sup>.

Notably, reaction is mixed: one commenter argues “distillation is not an attack” <sup>13</sup>, another claims DeepSeek’s use was “qualitatively different” than MiniMax’s and that lumping them together may be unfair <sup>14</sup>.

## 2) Standard Intelligence’s FDM-1: computer action learning from raw internet video

*Why it matters:* If agents can learn UI actions directly from video at internet scale, the bottleneck shifts from human labeling to compute and data access—enabling longer-horizon, higher-precision computer use.

Standard Intelligence announced **FDM-1**, a foundation model for computer actions that learns from **video** (not screenshot datasets with human action labels) <sup>1516</sup>. The reported approach uses:

- An **inverse dynamics model** to infer the action between frames <sup>17</sup>
- A **video encoder** that compresses nearly **2 hours** of high-res footage into the space other models use for **1 minute** <sup>18</sup>
- Auto-labeling to reach **11M hours** of screen recordings after training on **40k hours** of labeled data (described as **550,000×** larger than the biggest open dataset) <sup>19</sup>

Demonstrations include constructing a gear in Blender, finding software bugs, and driving a real car through San Francisco using arrow keys <sup>20</sup>. One claim

---

<sup>4</sup> post by @LiorOnAI  
<sup>5</sup> post by @AnthropicAI  
<sup>6</sup> post by @LiorOnAI  
<sup>7</sup> post by @LiorOnAI  
<sup>8</sup> post by @LiorOnAI  
<sup>9</sup> post by @AnthropicAI  
<sup>10</sup> post by @LiorOnAI  
<sup>11</sup> post by @LiorOnAI  
<sup>12</sup> post by @LiorOnAI  
<sup>13</sup> post by @glennko  
<sup>14</sup> post by @andersonbcdefg  
<sup>15</sup> post by @si\_pbc  
<sup>16</sup> post by @LiorOnAI  
<sup>17</sup> post by @LiorOnAI  
<sup>18</sup> post by @LiorOnAI  
<sup>19</sup> post by @LiorOnAI  
<sup>20</sup> post by @si\_pbc

highlights the driving result after **<1 hour** of training footage <sup>21</sup>. A separate framing: this takes computer action learning from “data-constrained” to “compute-constrained” <sup>22</sup>.

### 3) OpenAI ships WebSockets for agents + upgrades real-time voice with gpt-realtime-1.5

*Why it matters:* Latency and tool-call overhead are becoming core constraints for agent products; reducing round trips can translate directly into faster, more reliable agentic workflows.

OpenAI introduced **WebSockets in the Responses API** for low-latency, long-running agents with heavy tool calling <sup>23</sup>. The mode keeps a persistent connection and sends only incremental inputs rather than resending full context each turn <sup>24,25</sup>. OpenAI says maintaining in-memory state can speed up runs with **20+ tool calls by 20%–40%** <sup>26</sup>.

Early third-party results highlighted in posts:

- **Cline** reported **~15% faster** on simple tasks and **~39% faster** on complex multi-file workflows (best cases **50%**) vs the standard API, noting a small handshake overhead that amortizes on heavy tool use <sup>27,28</sup>.
- **Cursor** said OpenAI models are now up to **30% faster** after upgrading users to WebSockets <sup>29</sup>.

On voice, OpenAI released **gpt-realtime-1.5** in the Realtime API with improved instruction following, tool calling, and multilingual accuracy <sup>30</sup>. OpenAI also reported internal eval lifts: **+5%** on Big Bench Audio, **+10.23%** on alphanumeric transcription, and **+7%** on instruction following <sup>31</sup>. Partners shared deployment signals, including Genspark’s phone-call alpha test reporting a **66% human connection rate** (up from 43.7%) and a reduced “problem case rate” (2.1% vs 4.2%) <sup>32</sup>.

### 4) Guide Labs releases Sterling-8B, positioning interpretability as a first-class model feature

*Why it matters:* Claims of built-in traceability and memorization suppression—if they hold up in practice—target two recurring barriers to high-stakes deploy-

---

<sup>21</sup> post by @LiorOnAI

<sup>22</sup> post by @Hangsiin

<sup>23</sup> post by @OpenAIDevs

<sup>24</sup> post by @OpenAIDevs

<sup>25</sup> post by @cline

<sup>26</sup> post by @OpenAIDevs

<sup>27</sup> post by @cline

<sup>28</sup> post by @cline

<sup>29</sup> post by @leerob

<sup>30</sup> post by @OpenAIDevs

<sup>31</sup> post by @OpenAIDevs

<sup>32</sup> post by @genspark\_ai

ment: understanding *why* outputs occur and controlling training-data leakage.

Guide Labs announced **Steering-8B**, described as the first and largest “large-scale inherently interpretable” language model <sup>33</sup>. It is presented as tracing each generated token back to **input context, training data, and human-understandable concepts** <sup>34</sup>. The team also claims it can **self-monitor memorized content and suppress it at inference time without retraining** <sup>35</sup>.

Release links: Guide Labs post, GitHub, and Hugging Face <sup>363738</sup>.

## 5) DeepSeek V4 signals intensify alongside claims of training on NVIDIA Blackwell despite U.S. export restrictions

*Why it matters:* Reports of cutting-edge GPU access despite bans, combined with ongoing model-copying allegations, sharpen the question of what policy levers (compute vs access vs outputs) can realistically constrain capability diffusion.

Reuters reporting (as shared in posts) cites a senior U.S. official saying DeepSeek’s new model—described as **imminent**—was trained using **NVIDIA Blackwell GPUs** despite the export ban <sup>3940</sup>. Separately, posts suggested DeepSeek V4’s release is imminent, with one claim pointing to a pre-release polish pattern (merging PRs) <sup>41</sup>.

## Research & Innovation

### Monitoring reasoning traces: when chain-of-thought (CoT) helps—and when it doesn’t

*Why it matters:* As CoT is used for oversight, the question isn’t just “does the model show its work,” but whether monitors can reliably extract the right signals.

A paper summary shared by DAIR.AI formalizes CoT “monitorability” using information theory: mutual information between CoT and output is **necessary but not sufficient** for effective monitoring <sup>42</sup>. It identifies two failure modes—**information gap** and **elicitation error**—and proposes two training approaches (oracle-based rewards for transparency, and a label-free conditional

---

<sup>33</sup> post by @guidelabsai

<sup>34</sup> post by @guidelabsai

<sup>35</sup> post by @guidelabsai

<sup>36</sup> post by @guidelabsai

<sup>37</sup> post by @guidelabsai

<sup>38</sup> post by @guidelabsai

<sup>39</sup> post by @AndrewCurran\_

<sup>40</sup> post by @niubi

<sup>41</sup> post by @intheworldofai

<sup>42</sup> post by @dair\_ai

mutual information objective) that improve monitor performance without degrading reasoning traces <sup>43</sup><sup>44</sup>. Paper link: <https://arxiv.org/abs/2602.18297> <sup>45</sup>.

### Synthetic reasoning structure (ByteDance): “semantic isomers” and Mole-Syn

*Why it matters:* Long CoT training can become unstable if the model learns incompatible “reasoning structures,” even when surface-level solutions look similar.

ByteDance research is summarized as treating strong long CoT as having a **molecular-like internal structure** with three behaviors: deep reasoning, self-reflection, and self-exploration <sup>46</sup><sup>47</sup><sup>48</sup><sup>49</sup>. The summary warns that simply copying reasoning traces can fail—mixing traces from different models can destabilize training due to incompatible structures (“semantic isomers”) <sup>50</sup>. A proposed method, **Mole-Syn**, extracts transition patterns (deep reasoning → reflection → exploration) and generates new structured synthetic data without verbatim copying <sup>51</sup>. Paper link: <https://arxiv.org/abs/2601.06002> <sup>52</sup>.

### Speech NER under real-world diversity: SF Streets benchmark + a small-sample fix

*Why it matters:* Navigation and emergency dispatch failures can be driven by named-entity transcription errors, especially across diverse linguistic backgrounds.

Together Research introduced **SF Streets**, a benchmark for named entity recognition in speech across **15 models** <sup>53</sup>. Reported metrics include **39%** average error rate on street names, **18%** lower accuracy for non-English speakers, and mis-transcriptions landing you **2.4 miles** off target <sup>54</sup>. A proposed fix—cross-lingual style transfer with **<1,000** synthetic samples—yielded a **60% relative improvement** on Whisper-Large <sup>55</sup>. SF Streets and US Streets datasets are said to be releasing publicly <sup>56</sup>.

---

<sup>43</sup> post by @dair\_ai

<sup>44</sup> post by @dair\_ai

<sup>45</sup> post by @dair\_ai

<sup>46</sup> post by @TheTuringPost

<sup>47</sup> post by @TheTuringPost

<sup>48</sup> post by @TheTuringPost

<sup>49</sup> post by @TheTuringPost

<sup>50</sup> post by @TheTuringPost

<sup>51</sup> post by @TheTuringPost

<sup>52</sup> post by @TheTuringPost

<sup>53</sup> post by @togethercompute

<sup>54</sup> post by @togethercompute

<sup>55</sup> post by @togethercompute

<sup>56</sup> post by @togethercompute

## Evaluations: OpenAI stops reporting SWE-bench Verified

*Why it matters:* If a benchmark is contaminated or broken, “leaderboard progress” can diverge from real capability—and distort model selection.

OpenAI says SWE-bench Verified is saturated due to test-design issues and contamination from public repositories, and recommends reporting **SWE-bench Pro** instead <sup>5758</sup>. A separate audit summary shared in posts claims that after reviewing **27.6%** of frequently failed tasks, at least **59.4%** had flawed tests that reject correct solutions <sup>59</sup>.

## Products & Launches

### Pip-installable vector search: Alibaba open-sources Zvec

*Why it matters:* Making vector search a library (not a server) lowers adoption friction for local RAG, edge retrieval, and offline-first apps.

Alibaba open-sourced **Zvec**, described as a vector database you can **pip install** with no servers or Docker <sup>60</sup>. Performance claims shared include **8,000+ QPS** on **10M vectors** and “2×” the previous leader on VectorDBBench <sup>6162</sup>. Repo: <https://github.com/alibaba/zvec> <sup>63</sup>.

### LlamaIndex: LlamaAgents Builder adds file uploads for document workflows

*Why it matters:* Example documents as context can make “natural language workflow building” more grounded—especially for schema inference and validation.

LlamaIndex added **file upload support** to LlamaAgents Builder, letting users upload example docs so the agent can infer schema, validation rules, and pre/post-processing logic <sup>6465</sup>. The tool is positioned for scalable extraction with citations over complex documents, with user review before approval <sup>6667</sup>. Walkthrough + signup links were shared <sup>6869</sup>.

---

<sup>57</sup> post by @OpenAIDevs

<sup>58</sup> post by @OpenAIDevs

<sup>59</sup> post by @rasbt

<sup>60</sup> post by @LiorOnAI

<sup>61</sup> post by @LiorOnAI

<sup>62</sup> post by @LiorOnAI

<sup>63</sup> post by @LiorOnAI

<sup>64</sup> post by @llama\_index

<sup>65</sup> post by @jerryjliu0

<sup>66</sup> post by @jerryjliu0

<sup>67</sup> post by @jerryjliu0

<sup>68</sup> post by @jerryjliu0

<sup>69</sup> post by @llama\_index

## Image generation: Reve V1.5 reaches the top of Image Arena with 4K output

*Why it matters:* Arena performance combined with higher-resolution output can influence which models become default choices for commercial design workflows.

Reve launched **Reve V1.5**, a text-to-image model with output up to **4K resolution** <sup>70</sup>. It ranked **top 3** in Image Arena behind GPT-Image-1.5 and Nano Banana Pro variants <sup>71</sup>. Detailed scores: <https://arena.ai/leaderboard/text-to-image> <sup>72</sup>.

## Developer tooling highlights

- **Devin Review:** an AI-powered interface for understanding complex PRs; now supports fixing PRs inline by asking Devin to propose changes and applying them with one click <sup>73</sup>. Try: <http://devinreview.com> <sup>74</sup>.
- **LangSmith:** shipped native tracing for **Google ADK** agents <sup>75</sup>. Docs: <https://docs.langchain.com/langsmith/trace-with-google-adk> <sup>76</sup>.
- **OpenRouter:** launched “Effective Pricing,” estimating average provider costs based on cache pricing and cache hit rates <sup>77</sup>.

## Industry Moves

### OpenAI expands enterprise deployment via “Frontier Alliances”

*Why it matters:* Enterprise adoption often hinges on integration and change management, not just model quality.

OpenAI announced **Frontier Alliances** with **BCG, McKinsey, Accenture, and Capgemini** to deploy “OpenAI Frontier” to enterprises globally <sup>7879</sup>. The partnerships emphasize strategy, workflow redesign, system integration, and change management <sup>80</sup>. Announcement: <https://openai.com/index/frontier-alliance-partners> <sup>81</sup>.

A separate report link shared on X says OpenAI is hiring **hundreds of AI consultants** to boost enterprise sales <sup>8283</sup>.

---

<sup>70</sup> post by @reve

<sup>71</sup> post by @arena

<sup>72</sup> post by @arena

<sup>73</sup> post by @cognition

<sup>74</sup> post by @cognition

<sup>75</sup> post by @LangChain

<sup>76</sup> post by @LangChain

<sup>77</sup> post by @OpenRouter

<sup>78</sup> post by @bradlightcap

<sup>79</sup> post by @kimmonismus

<sup>80</sup> post by @kimmonismus

<sup>81</sup> post by @bradlightcap

<sup>82</sup> post by @srimuppidi

<sup>83</sup> post by @srimuppidi

## Salesforce Ventures invests in Sakana AI

*Why it matters:* Enterprise-focused AI labs are positioning around vertical credibility and deployment readiness.

Sakana AI announced an investment from Salesforce Ventures and said Salesforce will evaluate integrating Sakana’s enterprise technology into Salesforce’s global platform offerings <sup>84</sup>.

## Anthropic hiring: interpretability research engineers

*Why it matters:* As models become more central to critical workflows, internal understanding of model behavior is treated as infrastructure work.

Anthropic’s interpretability team is hiring ~10 **research engineers** (no prior interpretability experience required), targeting seasoned ML infrastructure engineers <sup>85</sup>.

## Policy & Regulation

### Pentagon–Anthropic tensions over Claude safeguards

*Why it matters:* Government adoption of frontier models is colliding with limits on surveillance and autonomy—potentially shaping what “acceptable safeguards” look like in high-stakes deployments.

An Axios-sourced report shared on X says the Pentagon threatened to ban Claude from classified systems and planned a meeting with Anthropic CEO Dario Amodei involving an “ultimatum” <sup>8687</sup>. The same reporting says Claude is described as the **only AI model available** in the military’s classified systems and the **most capable** model for sensitive defense and intelligence work <sup>88</sup>. Anthropic is said to be willing to loosen restrictions, while still walling off mass surveillance of Americans and autonomous weapons that fire without human involvement <sup>89</sup>.

### Export controls vs capability diffusion

*Why it matters:* If frontier capability can be replicated via API outputs—or trained on restricted hardware anyway—policy focus may shift toward access, enforcement, and monitoring.

Posts summarizing Anthropic’s position argue that illicit distillation can remove safeguards and feed capabilities into military/intelligence/surveillance systems

---

<sup>84</sup> post by @SakanaAILabs

<sup>85</sup> post by @ch402

<sup>86</sup> post by @JenGriffinFNC

<sup>87</sup> post by @AndrewCurran\_

<sup>88</sup> post by @JenGriffinFNC

<sup>89</sup> post by @JenGriffinFNC

<sup>90</sup>. Separately, Reuters reporting (as cited in posts) claims DeepSeek trained on NVIDIA Blackwell GPUs despite U.S. export restrictions <sup>91</sup><sup>92</sup>.

## Quick Takes

- **Anthropic AI Fluency Index:** tracked **11 behaviors** across thousands of Claude.ai conversations; one finding shared says **85.7%** of conversations exhibited iteration and refinement <sup>93</sup><sup>94</sup>.
- **Claude memorization extraction claim:** researchers reported extracting **95.8%** of *Harry Potter and the Sorcerer’s Stone* from Claude Sonnet <sup>95</sup>.
- **IBM volatility tied to AI modernization headlines:** a post said IBM stock fell **>10%** after claims that Claude can streamline COBOL code <sup>96</sup>.
- **Gemini training for educators:** Google said it’s making Gemini training available to **6 million** U.S. K–12 teachers and higher-ed faculty, with modular training and badges <sup>97</sup>.
- **Veo 3.1 templates:** Google said templates are rolling out in the Gemini app to provide a “visual foundation” for video creation <sup>98</sup><sup>99</sup>.
- **Qdrant 1.17:** announced features include relevance feedback queries, lower latency under heavy writes, and a cluster-wide telemetry API <sup>100</sup>.
- **Wispr Flow:** launched on Android as an AI voice dictation app; Richard Socher said he’s a “very happy user” <sup>101</sup><sup>102</sup>.

---

## Sources

1. post by @AnthropicAI
2. post by @niubi
3. post by @Bayesian0\_0
4. post by @LiorOnAI
5. post by @AnthropicAI
6. post by @glennko
7. post by @andersonbcdefg

---

<sup>90</sup> post by @AnthropicAI  
<sup>91</sup> post by @AndrewCurran\_  
<sup>92</sup> post by @niubi  
<sup>93</sup> post by @AnthropicAI  
<sup>94</sup> post by @dejavucoder  
<sup>95</sup> post by @ivanfioravanti  
<sup>96</sup> post by @KobeissiLetter  
<sup>97</sup> post by @Google  
<sup>98</sup> post by @Google  
<sup>99</sup> post by @GeminiApp  
<sup>100</sup> post by @qdrant\_engine  
<sup>101</sup> post by @tankots  
<sup>102</sup> post by @RichardSocher

8. post by @si\_pbc
9. post by @LiorOnAI
10. post by @Hangsiin
11. post by @OpenAIDevs
12. post by @OpenAIDevs
13. post by @cline
14. post by @leerob
15. post by @OpenAIDevs
16. post by @OpenAIDevs
17. post by @genspark\_ai
18. post by @guidelabsai
19. post by @AndrewCurran\_
20. post by @niubi
21. post by @intheworldofai
22. post by @dair\_ai
23. post by @TheTuringPost
24. post by @TheTuringPost
25. post by @togethercompute
26. post by @OpenAIDevs
27. post by @rasbt
28. post by @LiorOnAI
29. post by @LiorOnAI
30. post by @llama\_index
31. post by @jerryjliu0
32. post by @reve
33. post by @arena
34. post by @arena
35. post by @cognition
36. post by @cognition
37. post by @LangChain
38. post by @OpenRouter
39. post by @bradlightcap
40. post by @kimmonismus
41. post by @srimuppidi
42. post by @SakanaAILabs
43. post by @ch402
44. post by @JenGriffinFNC
45. post by @AndrewCurran\_
46. post by @AnthropicAI
47. post by @dejavucoder
48. post by @ivanfioravanti
49. post by @KobeissiLetter
50. post by @Google
51. post by @Google
52. post by @GeminiApp
53. post by @qdrant\_engine

54. post by @tankots
55. post by @RichardSocher