

AI Moves Into Phones, Desktops, and Secure Workflows

AI News Digest

2026-05-13

AI Moves Into Phones, Desktops, and Secure Workflows

By AI News Digest • May 13, 2026

Google and OpenAI pushed AI closer to the interface, while Microsoft, SAP, and NVIDIA focused on security and governance for enterprise agents. The digest also covers a new speech-agent benchmark, Isomorphic Labs' \$2.1B raise, sparse-LLM efficiency gains, and Hugging Face's 1 million dataset milestone.

The clearest pattern

Today's strongest signal is that major AI companies are pushing models out of the chat box and into the places where people already work: phones, desktops, voice channels, and enterprise systems. In parallel, the infrastructure underneath those systems is getting more specialized around security, efficiency, funding, and data [1, 2, 3, 4, 5].

AI moves closer to the interface

Google pushes Gemini deeper into Android—and experiments with an AI pointer

Google said Gemini Intelligence will automate multi-step tasks across apps and Chrome, fill forms in one tap, turn spoken thoughts into polished text with Rambler, and build custom widgets [1]. The rollout starts this summer on the latest Samsung Galaxy and Google Pixel phones, then expands later this year to Android watches, cars, glasses, and laptops [6].

Separately, Google DeepMind showed experimental demos of an AI-enabled pointer that can understand what is under the cursor and respond to shorthand like “fix this” or “move that,” with examples spanning PDFs, tables, recipes, handwritten notes, and paused video frames [7, 8, 9, 10].

Why it matters: Google is positioning AI less as a separate assistant window and more as an interaction layer across devices and apps [1, 6, 11, 12].

OpenAI moves Codex from code generation into computer use

OpenAI showed Codex using local GUI apps via mouse movement, clicks, and typing, extending it from commands and files into everyday desktop software [2]. The demo emphasized parallel work across multiple apps with a separate cursor that does not interrupt the user, plus permissioning that limits Codex to only the apps a user explicitly allows [2].

OpenAI also said computer use can leverage accessibility-framework data to understand interfaces more accurately—including off-screen elements—and work with fast non-multimodal models like Spark; the feature is available on Mac now, with Windows coming soon [2].

Why it matters: OpenAI is folding computer use into its main model stack rather than treating it as a separate experimental agent [2].



Computer use in Codex (0:05)

Voice agents are improving, but the ceiling is still low

Artificial Analysis launched -Voice to measure speech-to-speech models on realistic customer service tasks across airline, retail, and telecom scenarios, including

tool use and noisy audio conditions [3]. It said even the strongest models today resolve only about half of these scenarios end-to-end [3].

In the first leaderboard, xAI’s Grok Voice Think Fast 1.0 led at 52.1% success, ahead of GPT-Realtime-2 (High) at 39.8% and Gemini 3.1 Flash Live Preview - High at 37.7% [3]. xAI also said Grok is already handling live Starlink phone operations autonomously at scale [13].

Why it matters: The new benchmark shows real progress in speech agents, but it also quantifies how much reliability work remains before voice systems can consistently close complex service loops [3].

Enterprise agents are getting a security and governance layer

Microsoft says 100+ specialized agents helped find exploitable bugs

Microsoft announced a multi-model agentic security system that combines more than 100 specialized agents across frontier and custom models to find exploitable bugs, and said it delivered top performance on the CyberGym benchmark [14]. The company added that the system helped find and fix 16 vulnerabilities ahead of Patch Tuesday and is now opening a private preview for customers [14].

Why it matters: Microsoft is framing agentic security as coordinated specialist systems tied to measurable outcomes, not just general-purpose assistants [14].

SAP and NVIDIA focus on runtime controls for specialized agents

SAP and NVIDIA said SAP is embedding NVIDIA OpenShell into SAP Business AI Platform as an open-source runtime for securely developing and deploying autonomous agents [4]. OpenShell provides isolated execution environments, policy enforcement at the filesystem and network layers, and infrastructure-level containment; it will act as the runtime security layer for SAP AI agents, including custom agents built in Joule Studio [4].

NVIDIA also said its NemoClaw blueprint will be available directly in Joule Studio for custom agents in areas like finance, procurement, supply chain, and manufacturing [4].

Why it matters: The emphasis here is on control layers for production deployment—runtime isolation, policy enforcement, and governance—rather than raw autonomy alone [4].

Capital, efficiency, and data are scaling with the applications

Isomorphic Labs raises \$2.1B for AI drug discovery

Isomorphic Labs said it raised \$2.1B in new funding to accelerate AI-driven drug discovery, building on work that began with AlphaFold and extending it into a mission to reimagine drug discovery [15].

“I’ve always believed the No.1 application of AI should be to improve human health.” [15]

Why it matters: This is a significant funding signal for AI in biology and drug discovery, not just for horizontal model development [15].

Sakana AI and NVIDIA report sparse LLM speedups on H100s

Sakana AI and NVIDIA introduced TwELL (Tile-wise ELLPACK) and custom CUDA kernels designed to make sparse transformer language models fit GPU execution better [16, 17]. The team says feedforward layers can exceed 95% sparsity with little to no downstream performance loss, translating into more than 20% faster training and inference on H100 GPUs, along with lower peak memory use and energy consumption [16, 17].

Why it matters: This turns a familiar efficiency idea—LLM sparsity—into reported wall-clock gains on current production hardware [17, 16].

Hugging Face crosses 1 million public datasets

Hugging Face said it has now passed 1,000,000 public datasets, with petabytes of data being downloaded, analyzed, and used for training by millions of builders every day [5]. It also said the dataset count doubled in the past eight months, versus four years to reach the first 500,000, and linked that acceleration to better AI agents that make it easier to build, share, and use custom datasets [5].

Clément Delangue argued that better data is becoming the next bottleneck for people who want to build AI themselves instead of relying on APIs [5].

Why it matters: The open ecosystem’s next constraint may be shifting from access to models toward access to usable data [5].

Sources

1. X post by @sundarpichai
2. Computer use in Codex
3. X post by @ArtificialAnlys
4. NVIDIA and SAP Bring Trust to Specialized Agents
5. X post by @ClementDelangue

6. X post by @sundarpichai
7. X post by @GoogleDeepMind
8. X post by @GoogleDeepMind
9. X post by @GoogleDeepMind
10. X post by @GoogleDeepMind
11. X post by @GoogleDeepMind
12. X post by @demishassabis
13. X post by @XFreeze
14. X post by @satyanadella
15. X post by @demishassabis
16. X post by @SakanaAILabs
17. X post by @hardmaru