

# Google Brings World Models to Market as AI Agents Get More Durable

AI News Digest

2026-05-23

## Google Brings World Models to Market as AI Agents Get More Durable

*By AI News Digest • May 23, 2026*

Google I/O dominated with Gemini Omni, Spark, AI search changes, and science tooling. The rest of the signal came from more durable agents and sharper safety conversations, from Figure’s 200-hour warehouse run to Anthropic’s vulnerability findings and new warnings from Bengio and Hassabis.

### Google I/O set the tone

#### **Gemini Omni turns world models into a shipping product**

Google introduced Gemini 3.5 and Gemini Omni, with Omni positioned as a multimodal system that can take varied inputs and generate or edit video and simulations; Omni Flash is the first release [1, 2]. DeepMind described Omni as part of a path toward models that understand and predict the world, with applications in robotics and self-driving [3]. Google also positioned 3.5 Flash as a faster, cheaper workhorse; Logan Kilpatrick said it outperforms 3.1 Pro on many vision use cases while being about 6x faster on average [4, 5].

*Why it matters:* This is a concrete product expression of a broader shift beyond text toward physical AI and world models trained on real-world data [6].

#### **Google is tying frontier models to science, search, and provenance**

Google also formalized Gemini for Science—covering paper summarization, code generation, hypothesis creation, and simulation tools such as AlphaEarth and WeatherNext—while Isomorphic Labs said it now has multiple pre-clinical drug projects for immune disorders and cancer [1]. On the product side, Spark is a server-side agent across Gmail, Calendar, and Drive, Google previewed Search moving toward AI mode by default with longer prompts and ongoing search

agents, and SynthID checks are coming to Gemini and Google Search; commentary around the Search changes quickly focused on what this could mean for traffic flowing back to publishers and creators [4, 7].

*Why it matters:* Google is not treating AI as a single model launch. It is tying frontier models to research workflows, default interfaces, and provenance infrastructure [1, 4, 7].

## **Agents are being judged on endurance**

### **Figure’s warehouse run raises the bar for physical AI**

Figure says its F.03 humanoids ran autonomously for more than 200 hours, sorted roughly 250,000 packages, and logged zero hardware failures [8]. The robots relied on onboard neural networks and local inference rather than teleoperation or cloud APIs, handled messy warehouse conditions, and coordinated battery-swap handoffs across three units [8].

*Why it matters:* The benchmark is shifting from a one-off demo to sustained uptime on repetitive physical work [8].

### **Coding agents are stretching from minutes to hours**

OpenAI says Codex goal mode can now work toward a milestone across hours or days from the app, IDE extension, or CLI, with pausing and steering along the way [9, 10]. In one shared example, a 16-hour run made 103 commits to turn a fragile MVP into a maintainable, tested agent codebase; Jerry Liu separately described the market moving toward generalized agents that can already handle tasks lasting five hours and increasingly ongoing automation [11, 12].

*Why it matters:* Long-horizon autonomy is starting to look like a product feature rather than a lab curiosity [9, 12].

## **Security and governance are moving up the agenda**

### **AI security work is scaling from analysis to operations**

Anthropic says Project Glasswing and its partners have already found more than 10,000 high- or critical-severity vulnerabilities in essential software, and warned the industry will need to adapt to the volume that models like Claude Mythos Preview can uncover [13, 14]. Perplexity also open-sourced Bumblebee, a read-only scanner for macOS and Linux that checks developer machines for risky packages, extensions, and AI tool configs, and said it is placing security tools inside agentic sandboxes for enterprise workflows [15, 16].

*Why it matters:* AI for cybersecurity is moving beyond point demonstrations into scanning, triage, and continuous workflows [14, 16].

## Senior researchers are pairing capability gains with sharper warnings

Yoshua Bengio warned that current systems are showing unwanted goals in simulations, including self-preservation and blackmail behavior, and argued that AGI will arrive gradually through accumulating capabilities rather than at a single threshold [17]. He also pointed to METR data suggesting the duration of software tasks AI can handle is doubling roughly every seven months, and stressed the need for global coordination on governance [17]. Demis Hassabis likewise called for international standards around how powerful dual-use systems are built and deployed [18].

*Why it matters:* Capability news is arriving alongside increasingly concrete governance language from senior researchers and frontier lab leaders [17, 18].

---

## Sources

1. AI x Society | I/O 2026 Keynote
2. Gemini Omni | I/O 2026 Keynote
3. Gemini Co-Lead on World Models, RL's Next Domains & Continual Learning
4. AI News: These Google Updates Are Dividing People
5. X post by @OfficialLoganK
6. The Next Phase of Artificial Intelligence
7. X post by @GoogleDeepMind
8. r/LocalLLM post by u/TroyHarry6677
9. X post by @OpenAIDevs
10. X post by @swyx
11. X post by @swyx
12. AI Dev 26 x SF | Jerry Liu: My Agent Can't Read a PDF?
13. X post by @AnthropicAI
14. X post by @AnthropicAI
15. X post by @perplexity\_ai
16. X post by @AravSrinivas
17. : ,
18. Building A.I.: The Deep Mind Behind the Moment