

# Anthropic's Science Push, LongCat's Debut, and a July Model Wave

AI High Signal Digest

2026-07-04

## Anthropic's Science Push, LongCat's Debut, and a July Model Wave

*By AI High Signal Digest • July 4, 2026*

Anthropic's science push and LongCat-2.0's open frontier-scale debut led the day, while multiple signs pointed to an imminent GPT-5.6 and Gemini release window. The brief also covers Leanstral, HOLA, ARIA, and changing economics around open-model deployment.

### Top Stories

*Why it matters: frontier competition is widening from chatbots to science workflows, domestic compute stacks, and back-to-back flagship launches.*

- **Anthropic is pushing deeper into life sciences.** At its AI for Science event, it introduced **Claude Science**, a workbench that brings fragmented scientific tools and datasets into one environment and can generate figures and visuals. Separate posts said Anthropic is moving from selling AI tools to drugmakers toward developing drugs itself, including plans to discover treatments for neglected diseases [1].
- **LongCat-2.0 pairs frontier-scale specs with an open release.** The model is a 1.6T-parameter MoE with ~48B active params, 1M context, sparse attention, and specialized expert groups for agentic coding. Reported scores included 70.8 on Terminal-Bench 2.1 and 59.5 on SWE-bench Pro, and the weights were released on Hugging Face; one post described the run as the largest known pretraining on non-Western chips [2, 3].
- **A July flagship-model cycle is taking shape.** Posts this week said OpenAI is targeting **GPT-5.6** for July 7-9 with more generous plan limits and stronger safeguards, while DeepMind tentatively moved **Gemini 3.5 Pro** to July 17 after a new pretrain. Sol, Terra, and Luna labels are also already appearing in Codex code and in at least one user UI screenshot

[4, 5, 6].

## Research & Innovation

*Why it matters: the strongest technical work today targeted formal reasoning, memory, and training reliability rather than just bigger scale.*

- **Mistral released Leanstral 1.5 for theorem proving.** The Apache-2 6B-active model hit SoTA on FATE-H/X, solved 587 of 672 PutnamBench problems at 10x lower cost, saturated miniF2F, and shipped with an open tech report plus LeanstralSafeVerify and FLTEval [7, 8].
- **HOLA offers a cleaner memory fix for linear attention.** The architecture combines a compressive recurrent state with a bounded exact KV cache, improving long-range recall while keeping O(1) efficiency. At 340M parameters, it lowered Wikitext perplexity to 22.92 from 27.32 and stayed robust on 32k-token RULER recall [9].
- **An MIT-style RLVR fix aims to preserve quality without reward hacking.** The approach adds an adversarial discriminator trained on human demonstrations, so the generator optimizes both task accuracy and human-likeness. Reported results across bug fixing, story generation, and a reward-hacking benchmark preserved accuracy gains while restoring diversity and reducing misbehavior [10].

## Products & Launches

*Why it matters: the most interesting product launches are starting to act on live context and feedback, not just generate outputs.*

- **CoreWeave ARIA** launched as an AI research agent inside W&B dashboards. It reads prior runs, identifies what is working, proposes the next experiments, and launches them; a demo showed parallel ARIA instances running on Karpathy’s nanochat and evaluating validation loss [11].
- **Poolside’s Laguna M.1** is now available free through Cline. The model has 225B parameters, 256k context, and is positioned for agentic coding and long-horizon work; Cline remains open source and bring-your-own-key [12, 13].
- **Opik** is pitching a more practical eval stack for production AI. The open-source tool covers observability for LLM apps, RAG systems, and agent workflows, and its Test Suites feature uses plain-English assertions with pass/fail outputs instead of raw scores [14, 15].

## Industry Moves

*Why it matters: deployment strategy and model economics are becoming as important as raw capability.*

- **Open models are gaining share, but usage is getting more disciplined.** One post said open-model usage rose from 10% to 30% of

AI tokens in a year, while another said companies are moving out of a token-maxing phase and toward cost controls plus task-specific portfolios of smaller and open models [16, 17].

- **Tencent Cloud will serve DeepSeek-V4 directly from DeepSeek’s own network starting mid-July.** Separate commentary read that as a sign DeepSeek’s compute cluster is expanding beyond earlier third-party GPU dependence [18, 19].
- **Cohere is leaning into customer-side deployment as a security differentiator.** The company said it deploys models directly to customers instead of having customers send data back to Cohere [20].

## Quick Takes

*Why it matters: these smaller updates still show where latency, safety practice, and open-model usability are heading.*

- **Qwen3-Omni serving got much faster:** replicating only the speech stages under load cut first-audio latency from about 6 seconds to about 0.6 seconds and lifted throughput about 5.4x on the same GPUs [21].
- **GLM-5.2 is now selectable in Claude Code via Hugging Face providers,** and one user said they now use it almost daily and have moved completely to open models [22, 23].
- **Safety reporting remains a strong lab norm:** one post cited a 319-page Fable model card, a 77-page GPT-5.6 model card, and a 26-page Gemini 3 safety report [24, 25].
- **DeepSeek V4 release names surfaced in the wild:** deepseek-v4-pro-202606 and deepseek-v4-flash-202605, alongside mention of a coding-focused plan [26].

---

## Sources

1. X post by @kimmonismus
2. X post by @Meituan\_LongCat
3. X post by @teortaxesTex
4. X post by @synthwavedd
5. X post by @testingcatalog
6. X post by @DevAdventur3s
7. X post by @AlbertQJiang
8. X post by @mertunsal2020
9. X post by @omarsar0
10. X post by @dair\_ai
11. X post by @wandb
12. X post by @cline
13. X post by @cline
14. X post by @dl\_weekly

15. X post by @dl\_weekly
16. X post by @togethercompute
17. X post by @MTSlive
18. X post by @tphuang
19. X post by @teortaxesTex
20. X post by @cohere
21. X post by @vllm\_project
22. X post by @zRdianjiao
23. X post by @\_akhaliq
24. X post by @andy\_l\_jones
25. X post by @idavidrein
26. X post by @teortaxesTex