

# Qwen3.5’s open-weight multimodal MoE lands across inference stacks as agent efficiency and inference economics dominate

AI High Signal Digest

2026-02-17

## Qwen3.5’s open-weight multimodal MoE lands across inference stacks as agent efficiency and inference economics dominate

*By AI High Signal Digest • February 17, 2026*

Qwen3.5’s open-weight multimodal release rapidly landed across vLLM/SGLang and major hosted endpoints, while MiniMax M2.5’s agent-focused positioning translated into broad distribution and usage signals. Meanwhile, inference economics (Blackwell Ultra claims, cross-vendor benchmarking, and memory constraints) and agent-efficiency research (trajectory pruning, adaptive reasoning, multi-step tool-use benchmarks) were recurring themes.

### Top Stories

#### 1) Qwen3.5’s open-weight release becomes a cross-stack default (model + inference tooling + hosted endpoints)

*Why it matters:* A strong open-weight model only reshapes the landscape when it’s easy to run (local + cloud) and supported by the major inference stacks.

Alibaba released **Qwen3.5-397B-A17B**, the first open-weight model in the Qwen 3.5 series, positioned as a **native multimodal** model trained for real-world agents and licensed under **Apache 2.0** [1]. Multiple posts emphasize the architecture: **hybrid linear attention + sparse MoE** with **397B total parameters / 17B active**, targeting throughput and latency [2].

Distribution and “day-0” enablement showed up quickly:

- **Inference stacks:** day-0 support in **vLLM** (with deployment recipes) [2] and **SGLang** (with cookbook + PR) [3].

- **Hardware ecosystems:** AMD published day-0 guidance for **Instinct GPUs** via SGLang/vLLM [4]. NVIDIA highlighted a free build surface and a NeMo fine-tune config [5].
- **Hosted endpoints:** Together AI listed it as production-ready with a 99.9% SLA and highlighted **87.8% MMLU-Pro** plus “early fusion” multimodality [6]. OpenRouter also added Qwen3.5-397B-A17B and Qwen3.5 Plus variants [7].
- **Local options:** Unsloth published local run artifacts (including a GGUF) and claimed 4-bit can run on **256GB RAM** [8]. Ollama made it available on cloud with `ollama run qwen3.5:cloud` [9].

Context + efficiency were recurring themes: Qwen3.5 Plus is described as the hosted API variant of the 397B model with **1M context** (vs native **256K**) plus search and code interpreter [10]. Separately, an analysis of Qwen3.5’s **KV-cache** footprint at **262K** context estimates **8.05 GB** in BF16 (or **4.025 GB** in FP8 KV) and attributes it to the use of **45 gated deltanet layers** [11].

Pricing sparked debate. Novita listed **\$0.6/M input** and **\$3.6/M output** tokens [12], while another post noted a **China vs international** API price delta (e.g., **\$0.12 prefill + \$0.69 decode** in China vs **\$0.4/\$2.4** internationally) [13].

## 2) “Open model week” continues: MiniMax M2.5’s agentic coding positioning turns into usage share

*Why it matters:* Popularity signals (usage, partners, integrations) often matter as much as model evals when developers choose defaults.

MiniMax **M2.5** is repeatedly framed as strong for agent workflows. MiniMax says it uses **per-token process rewards** to better utilize signal across reasoning steps [14], and a separate post claims “frontier coding performance” at **at least 1/10th the cost** of closed-source models [15].

Adoption signals:

- OpenRouter reported M2.5 became the **most popular model** and hit **#1** on the weekly leaderboard in **four days** [16, 17].
- Together AI announced MiniMax M2.5 availability for production-scale agentic workflows [18, 19].
- Baseten listed M2.5 on its model APIs [20].
- Windsurf added **GLM-5** and **MiniMax M2.5** with limited-time credit discounts [21].

## 3) Inference economics take center stage: Blackwell Ultra claims, cross-vendor benchmarking, and memory constraints

*Why it matters:* For agents and long-context workloads, cost-per-token and performance-per-watt increasingly set the ceiling for deployment.

NVIDIA highlighted **Blackwell Ultra GB300 NVL72**, claiming up to **50× higher performance per megawatt** (also framed as “tokens per watt”) and **35× lower cost per token** versus Hopper, aimed at low-latency and long-context agentic use cases [22, 23, 24].

In parallel, SemiAnalysis promoted **InferenceX v2** benchmarks comparing **Blackwell vs AMD vs Hopper**, covering systems like **GB300 NVL72, MI355X, B200, H100**, and techniques such as disaggregated serving and wide expert parallelism, tested across **SGLang, vLLM, TRTLLM** [25].

Supply-side constraints also surfaced:

- Western Digital reportedly sold out its entire **2026 hard drive capacity**, with most supply locked by top enterprise customers (consumers ~5% of revenue) [26].
- A Bloomberg-linked thread described a growing **memory chip crisis**, with Sony considering pushing PS6 to **2028 or 2029** [27].

#### 4) Agent reliability and efficiency become the research focal point (benchmarks + pruning + adaptive reasoning)

*Why it matters:* As agents move into longer-horizon workflows, the bottleneck shifts to multi-step execution quality and wasted tool calls/tokens.

Notable signals this cycle:

- **WebClipper** models web-agent search as state graphs and prunes into minimal DAGs, reporting **~20% reduction in tool-call rounds** while maintaining or improving accuracy; it also introduces **F-AE Score** to balance accuracy and efficiency in trajectories [28].
- **CogRouter** dynamically adjusts reasoning depth step-by-step across four cognitive levels; the report cites a **7B** model reaching **82.3%** success on agent benchmarks while using **62% fewer tokens** than a baseline it outperformed (GPT-4o is named in the post) [29].
- **SciAgentGym** evaluates multi-step scientific tool use with **1,780 tools** across **4 disciplines**; a post claims success drops from **60.6% to 30.9%** as interaction steps increase, and presents **SciForge** (dependency-graph trajectory synthesis) with an **SciAgent-8B** outperforming much larger models on scientific workflows [30].

#### 5) OpenAI adds an enterprise security posture for ChatGPT workflows

*Why it matters:* Prompt-injection and connected-app risks are increasingly operational issues; “security modes” change what’s viable for regulated deployments.

OpenAI introduced **Lockdown Mode** for ChatGPT (initially for enterprise/business users), disabling some tools/capabilities that could be exploited

to exfiltrate sensitive data via attacks such as prompt injection, including switching to cached browsing and limiting broader web interaction [31].

---

## Research & Innovation

*Why it matters:* The highest leverage work right now targets faster inference (one-step generation, sparse routing), better reasoning under constraints (few-parameter adaptation, RL objectives), and agent training that wastes fewer steps.

### Highlights from recent papers (curated list)

- **Generative Modeling via Drifting:** a generative framework that evolves the pushforward distribution to enable **native one-step inference**, reporting **FID 1.54 (latent) / 1.61 (pixel)** on ImageNet  $256 \times 256$  [32].
- **TinyLoRA (“Learning to Reason in 13 Parameters”):** scales low-rank adapters down to a single parameter; the post claims **91% GSM8K** accuracy with **13 trained parameters** and recovery of **90%** of gains on AIME/MATH500 while training  $1000 \times$  fewer parameters than typical LoRA approaches [33].
- **Maximum Likelihood Reinforcement Learning:** defines compute-indexed objectives interpolating RL and exact maximum likelihood; claims up to  **$20 \times$  test-time scaling efficiency gains** over GRPO in math reasoning and code generation tasks [34].
- **Kimi K2.5 (agentic multimodal training):** combines text/vision training stages and a parallel “Agent Swarm” orchestration approach; claims  **$4.5 \times$  latency reduction** vs single-agent baselines [35].
- **Generative meta-models of LLM activations:** trains diffusion models on **1B residual stream activations** to learn priors over internal states; reported to improve fluency for steering interventions and to scale sparse probing as neurons isolate concepts [36].
- **On-Policy Context Distillation (OPCD):** trains on self-generated trajectories while minimizing reverse KL vs a context-conditioned teacher; claims improved math reasoning/text-games accuracy and stronger OOD performance, supporting cross-size distillation [37].
- **SkillRL:** builds a hierarchical skill library (SkillBank) from trajectories using retrieval + recursive evolution; claims **+15.3%** on ALFWorld/WebShop and reduced token footprint [38].
- **Retrieval-aware distillation for Transformer–SSM hybrids:** keeps only attention heads critical for in-context retrieval; claims retaining **2%** of heads recovers **95%+** of teacher performance, enabling **5–6 $\times$**  memory efficiency in retrieval-heavy settings [39].
- **ViT-5:** modernizes ViT with changes to normalization/activation/gating while preserving Attention–FFN layout; claims **84.2%** ImageNet-1k top-1

accuracy and **FID 1.84** in an SiT diffusion framework [40].

### Additional notable research signals

- **Deep-Thinking Ratio (DTR)** and **Think@n**: proposes a measure of “deep thinking” effort and a test-time strategy that prefers/aggregates higher-DTR generations while stopping early on low-DTR ones; claims to outperform self-consistency using **~50% less compute** [41, 42].
  - **MonoLoss**: a plug-in objective for SAEs that rewards semantically consistent activations to increase monosemanticity (MonoScore) across SAEs trained on CLIP/SigLIP2/ViT features [43].
  - **Skills as procedural knowledge**: a benchmark across **86 tasks / 11 domains** and **7,300+ trajectories** reports curated skills improve pass rates by **16.2pp** on average, while self-generated skills show no average benefit; concise skills outperform comprehensive docs [44].
- 

## Products & Launches

*Why it matters:* Agents and models become “real” when they’re packaged into deployable workflows: chat surfaces, developer tooling, hosted inference, and observability.

### Agents in mainstream chat surfaces

- **Manus Agents** launched inside chat apps (starting with Telegram), offering long-term memory, multi-step execution, and tool integrations (Gmail, Calendar, Notion, etc.) [45, 46]. One post also claimed this helps explain **Meta’s acquisition of Manus** (attributed as commentary) [46].
- **SkyBot (Skywork)** is positioned as a cloud-native agent for long-term task execution with **zero setup** (no code/keys/servers), running in the background and accessible via phone/Discord/Telegram [47, 48].

### Voice and realtime interaction tooling

- **NVIDIA PersonaPlex** is live on fal as a full-duplex model that listens and speaks simultaneously, handles interruptions/backchannels, and aims for low-latency spoken interaction with a consistent persona [49].

### Agent harness + “background agent” infrastructure

- **Ollama** added **subagents and web search** in Claude Code, enabling parallel tasks in isolated contexts (file search, code exploration, research) and automatic web search via the Anthropic compatibility layer—no MCP servers or API keys required [50, 51].

- **Terminal Use** (YC launch) provides infrastructure for background agents, including filesystem forking and parallel agent runs; YC’s framing emphasizes that agent apps often “win on the harness” rather than the model alone [52].

### Developer tools and document workflows

- **LlamaIndex / LlamaCloud** released a parsing feature to convert complex PDFs (tables, charts, multi-column layouts) into clean markdown/JSON via clickable templates [53].
- **Base44** launched a standalone backend (CLI-first, realtime, AI-agent friendly) for deploying auth/database/hosting from the CLI [54].
- **Synthesia** launched a Word-to-video flow (upload → adjust settings → generate video; optional brand kit/translation/interactivity) [55].

### Industry Moves

*Why it matters:* Usage momentum, distribution partnerships, and supply constraints influence which models and tools become defaults.

- **Codex usage momentum:** OpenAI’s Sam Altman said Codex weekly users have **more than tripled** since the beginning of the year [56].
- **Anthropic expansion:** Anthropic opened a **Bengaluru office**, calling India its second-largest Claude.ai market [57].
- **Model availability as a differentiator:** Artificial Analysis benchmarked **Kimi K2.5** across **8 providers**, showing output speed variation of **~330 tokens/s**, plus differences in latency (TTFAT/TTFT), pricing, context support, and multimodality/tooling support [58].
- **Autonomy in the physical world:** Waymo began full autonomous operations with its **6th gen platform**; one post claims Waymo does **500,000+ driverless rides/week**, and cites a **~\$70k** per-vehicle cost with room to fall by ~50% over two years [59, 60].

### Policy & Regulation

*Why it matters:* As deployments expand, governance shows up as procurement rules, contract terms, transparency norms, and security controls.

- **Pentagon–Anthropic tension (Axios-linked):** posts report the Pentagon said Anthropic will “pay a price,” while Anthropic is described as willing to loosen terms of use but seeking safeguards against mass surveillance of Americans and fully autonomous weapons [61, 62]. Another Axios-linked post says the Pentagon is considering labeling Anthropic a “supply chain risk,” which could force vendors to cut ties [63].

- **Peer review integrity:** a post claims ICML journal editors inserted **hidden prompt injections** into papers to detect AI-assisted reviewing, causing at least one reviewer to consider desk rejection after discovering it [64].
- **Labor-market signal:** one post claims the US BLS revised 2025 job numbers downward by **over 1 million**, with the Information sector revised down **88,000 jobs (3%)**, attributed by economists in the post to AI automation of tech-heavy roles [65].

---

## Quick Takes

*Why it matters:* Smaller signals often become the next “normal”—or the next operational risk.

- **ByteDance BitDance:** an open-source **autoregressive image model** with a GitHub repo and details like up to **32× downsampling** and a codebook size up to **2<sup>256</sup>** [66, 67].
- **Chinese humanoid robots (Spring Festival Gala):** Unitree showed a large robot cluster; posts highlight an autonomous Kung Fu performance and an H2 “Monkey King” segment with robot dogs [68].
- **Qwen3.5 in evaluation venues:** LM Arena added Qwen3.5-397B-A17B to Text/Vision/Code arenas and asked users to test and vote for leaderboards [69].
- **Anthropic research on skill retention:** a post summarized an RCT where AI coding assistance decreased skill mastery by **17%** among 52 software engineers, with debugging most affected despite minimal productivity gains [70].
- **Meta patent:** a post claims Meta patented an AI that takes over a deceased person’s account to keep posting and chatting [71].

---

## Sources

1. X post by @Alibaba\_Qwen
2. X post by @vllm\_project
3. X post by @lmsysorg
4. X post by @AIatAMD
5. X post by @NVIDIAAIDev
6. X post by @togethercompute
7. X post by @OpenRouterAI
8. X post by @UnsllothAI
9. X post by @ollama
10. X post by @JustinLin610
11. X post by @bnjmn\_marie
12. X post by @novita\_labs

13. X post by @sun\_hanchi
14. X post by @MiniMax\_AI
15. X post by @basetenco
16. X post by @MiniMax\_AI
17. X post by @OpenRouterAI
18. X post by @togethercompute
19. X post by @MiniMax\_AI
20. X post by @basetenco
21. X post by @windsurf
22. X post by @kimmonismus
23. X post by @nvidia
24. X post by @kimmonismus
25. X post by @SemiAnalysis\_
26. X post by @kimmonismus
27. X post by @tomwarren
28. X post by @dair\_ai
29. X post by @omarsar0
30. X post by @dair\_ai
31. X post by @cryps1s
32. X post by @TheAITimeline
33. X post by @TheAITimeline
34. X post by @TheAITimeline
35. X post by @TheAITimeline
36. X post by @TheAITimeline
37. X post by @TheAITimeline
38. X post by @TheAITimeline
39. X post by @TheAITimeline
40. X post by @TheAITimeline
41. X post by @WeiLin\_\_Chen
42. X post by @WeiLin\_\_Chen
43. X post by @iScienceLuvr
44. X post by @omarsar0
45. X post by @ManusAI
46. X post by @kimmonismus
47. X post by @kimmonismus
48. X post by @kimmonismus
49. X post by @fal
50. X post by @ollama
51. X post by @ollama
52. X post by @ycombinator
53. X post by @jerryjliu0
54. X post by @MS\_BASE44
55. X post by @synthesiaIO
56. X post by @sama
57. X post by @AnthropicAI
58. X post by @ArtificialAnlys

59. X post by @reed
60. X post by @fchollet
61. X post by @unusual\_whales
62. X post by @kimmonismus
63. X post by @DavidLawler10
64. X post by @paul\_cal
65. X post by @kimmonismus
66. X post by @aisearchio
67. X post by @teortaxesTex
68. X post by @XRoboHub
69. X post by @arena
70. X post by @dl\_weekly
71. X post by @kimmonismus