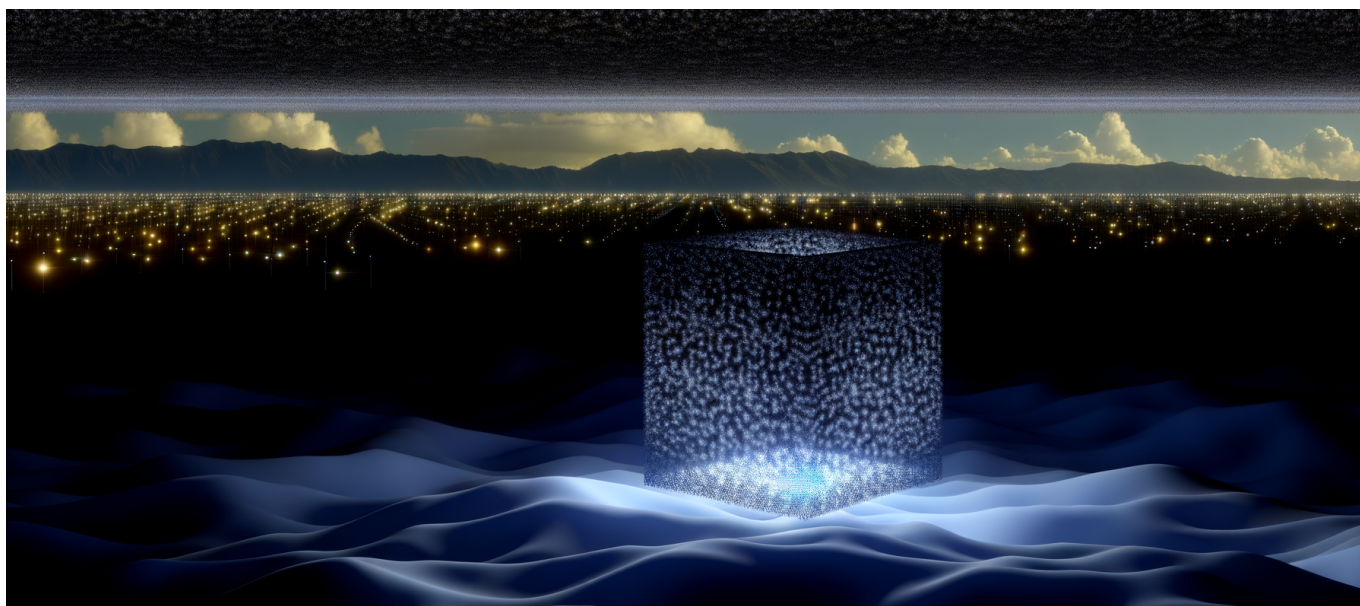




Liquid AI Releases LFM2.5-1.2B-Thinking: 1.2B Parameter Reasoning Model Runs Fully Offline in 900 MB on Phones



Liquid AI Releases LFM2.5-1.2B-Thinking: 1.2B Parameter Reasoning Model Runs Fully Offline in 900 MB on Phones

A reasoning model smaller than a mobile game just outscored models twice its size—and it runs entirely on your phone with no internet connection. Liquid AI's new release marks the moment edge AI stopped being a compromise.

The News: Sub-1GB Reasoning Arrives

[Liquid AI released LFM2.5-1.2B-Thinking on January 20, 2026](#), an open-source reasoning model with 1.17 billion parameters that fits in approximately 900 MB of memory. The model scores 87.96 on MATH 500 and 85.60 on GSM8K—benchmarks that measure mathematical reasoning and grade-school math problem-solving respectively.



Liquid AI Releases LFM2.5-1.2B-Thinking: 1.2B Parameter Reasoning Model Runs Fully Offline in 900 MB on Phones

These numbers matter because they exceed the performance of models with significantly more parameters. [According to Liquid AI's benchmarks](#), LFM2.5-1.2B-Thinking matches or exceeds Qwen3-1.7B in thinking mode—a model with 40% more parameters. It outperforms Granite-4.0-1B, Gemma-3-1B-IT, and Llama-3.2-1B Instruct on reasoning tasks.

The model runs at 239 tokens per second on AMD CPUs and 82 tokens per second on mobile NPUs. It supports a 32,768 token context window across eight languages. Open weights are available on Hugging Face with inference support for llama.cpp, MLX, vLLM, and hardware from AMD, Qualcomm, Apple, and Nvidia.

Why This Matters: The Economics of Edge AI Just Changed

The conventional wisdom in AI deployment has been simple: serious reasoning requires serious infrastructure. Cloud GPUs for inference, network latency as an accepted cost, and API pricing models that scale with usage. LFM2.5-1.2B-Thinking directly challenges each of these assumptions.

The Cost Equation

Consider what it costs to serve reasoning capabilities today. A single A100 GPU for inference runs \$1-3 per hour in cloud environments. For applications requiring thousands of concurrent users, costs escalate rapidly. Each API call carries latency—typically 500-2000ms for round trips to cloud inference endpoints.

A model running locally at 82 tokens per second on a mobile NPU eliminates these costs entirely. The inference hardware is already paid for—it's the phone in your users' pockets. The marginal cost of each additional reasoning call drops to the electricity cost of running the NPU, which is measured in fractions of cents.

Who Wins

Mobile developers can now build applications that think without network dependencies. Consider medical apps in rural areas, educational tools in regions with unreliable connectivity, or industrial applications in environments where cloud connectivity is impossible or prohibited.



Privacy-focused enterprises gain the ability to deploy reasoning capabilities without data ever leaving the device. Financial institutions, healthcare providers, and legal firms that previously couldn't use cloud AI due to compliance constraints now have an on-ramp.

Hardware manufacturers finally have a compelling reason for users to care about NPU specifications. The NPU has been the orphan silicon in mobile chips—present but underutilized. A model like this gives that hardware a purpose consumers will actually notice.

Who Loses

API-first AI companies face margin compression. If edge models can handle routine reasoning tasks, cloud APIs become reserved for only the most complex queries. The volume of billable API calls drops.

The “bigger is better” narrative takes a hit. The industry has spent years conditioning customers to believe that parameter count correlates with capability. A 1.2B model outperforming models at 40% more parameters complicates that marketing story.

Technical Depth: How LFM2.5-1.2B-Thinking Achieves Its Efficiency

The performance of this model isn't accidental. It emerges from specific architectural decisions that prioritize efficiency without sacrificing reasoning capability.

Architecture Design

[The model uses a 16-layer architecture](#) combining two distinct block types: 10 double-gated LIV convolution blocks and 6 grouped-query attention (GQA) blocks. This hybrid approach is unusual—most reasoning models rely entirely on attention mechanisms.

LIV convolutions (Linear-Invariant Variational convolutions) are Liquid AI's proprietary contribution. They provide sequence modeling capabilities similar to attention but with computational characteristics that scale more favorably for long



sequences. By handling most of the sequence processing through convolutions and reserving attention for specific layers, the model reduces the quadratic scaling problem that makes large context windows expensive.

The GQA blocks use grouped-query attention, a technique that reduces the memory footprint of attention mechanisms by sharing key-value heads across query heads. This is the same approach used in Llama 2 and other efficient transformer variants.

The 65,536 token vocabulary provides adequate coverage across eight languages while keeping embedding matrices manageable. Training used 28 trillion tokens—a substantial dataset that explains the model’s performance density.

The RLVR Training Innovation

Perhaps the most significant technical achievement is the doom loop reduction. **Doom loops—where a model gets stuck generating repetitive or circular outputs—dropped from 15.74% to 0.36% through RLVR (Reinforcement Learning from Verifiable Rewards) training.**

This matters because reasoning models are particularly susceptible to doom loops. The extended chain-of-thought reasoning that makes these models useful also creates more opportunities for the model to enter degenerate states. A doom loop rate under 0.5% means the model can be deployed in production settings without requiring extensive output monitoring and regeneration logic.

RLVR differs from RLHF (Reinforcement Learning from Human Feedback) in that it uses automatically verifiable rewards rather than human preference data. For mathematical reasoning, correctness is often verifiable programmatically. The model can be rewarded for reaching correct answers through valid reasoning steps, creating a training signal that’s both scalable and precise.

Benchmark Context

The 87.96 MATH 500 score deserves context. [For comparison](#), the non-thinking variant (LFM2.5-1.2B-Instruct) scores approximately 63 on the same benchmark. That 25-point improvement from enabling thinking mode represents the value of test-time compute—letting the model reason through problems rather than answering immediately.



Liquid AI Releases LFM2.5-1.2B-Thinking: 1.2B Parameter Reasoning Model Runs Fully Offline in 900 MB on Phones

The 85.60 GSM8K score places this model in competitive range with much larger systems from even a year ago. GSM8K specifically tests grade-school math word problems, requiring the model to parse natural language, identify relevant quantities, and apply correct operations.

Instruction following (Multi IF score of 69) and tool use (BFCLv3 score of 57) indicate the model handles structured interactions competently, though these aren't its primary strengths. The tool use score in particular suggests room for improvement in function-calling scenarios.

The Contrarian Take: What the Coverage Gets Wrong

The immediate response to this release has focused on the “small model, big performance” narrative. That’s the easy story. It’s also incomplete.

What’s Overhyped: The Benchmark Numbers

MATH 500 and GSM8K are useful benchmarks, but they test a narrow slice of reasoning capability. Mathematical reasoning is particularly amenable to the chain-of-thought approach because math problems have verifiable correct answers and clear logical steps.

Real-world reasoning often lacks these properties. Ambiguous premises, incomplete information, and problems where “correct” depends on context—these challenge models in ways that mathematical benchmarks don’t capture. A CTO evaluating this model for business applications should test extensively on their actual use cases, not assume benchmark performance translates directly.

The 40% parameter efficiency claim also requires scrutiny. Qwen3-1.7B uses a different architecture with different design tradeoffs. Comparing parameter counts across architecturally distinct models is like comparing engine displacement across vehicles with different transmission designs—it’s not meaningless, but it’s not the whole story either.

What’s Underhyped: The Inference Speed

The real breakthrough is 239 tokens per second on AMD CPUs. This number gets



Liquid AI Releases LFM2.5-1.2B-Thinking: 1.2B Parameter Reasoning Model Runs Fully Offline in 900 MB on Phones

buried in coverage focused on benchmark scores, but it's arguably more consequential.

Most edge AI discussions assume you need specialized hardware—NPUs, TPUs, or at minimum a capable GPU. A model that runs at 239 tokens per second on commodity CPUs can deploy on virtually any hardware manufactured in the last five years. This includes older laptops, budget Android devices, industrial controllers, and embedded systems.

82 tokens per second on mobile NPUs is respectable but expected. 239 tokens per second on CPUs is exceptional and opens deployment scenarios that dedicated AI hardware doesn't reach.

What's Actually New: The Architecture Bet

Liquid AI's hybrid convolution-attention architecture represents a genuine technical divergence from the transformer orthodoxy. Most efficiency improvements in recent years have come from training techniques, quantization, and attention optimizations. Liquid AI is betting that the fundamental architecture can be improved.

This bet matters because it's falsifiable. If LIV convolutions provide genuine advantages for on-device inference, we should see other labs adopting similar approaches. If they don't, and pure transformer architectures continue to improve faster, Liquid AI's architectural differentiation becomes a liability rather than an advantage.

The next 12-18 months will tell us whether this is a meaningful innovation or an interesting dead end.

Practical Implications: What You Should Do

For Engineers Building AI Features

Test this model against your actual workloads. Download the weights from Hugging Face and run inference on representative inputs from your domain. Pay attention to failure modes, not just success cases.

The llama.cpp integration means you can run inference without complex setup. A



Liquid AI Releases LFM2.5-1.2B-Thinking: 1.2B Parameter Reasoning Model Runs Fully Offline in 900 MB on Phones

basic evaluation pipeline:

1. Pull the model weights in GGUF format
2. Set up llama.cpp with your target hardware (CPU, Metal for Mac, or CUDA if available)
3. Create a test suite from real user queries or business scenarios
4. Measure latency, accuracy, and output quality
5. Compare against your current solution's cost and performance

If you're building mobile applications, the Qualcomm and Apple silicon support means you can deploy today without waiting for vendor-specific optimizations. Test on actual target devices, not just development hardware.

For CTOs Evaluating Architecture Decisions

The existence of capable edge models changes build-vs-buy calculations. Questions to ask:

What fraction of your AI inference could run on-device? Most applications have a bimodal distribution—some queries require full cloud capabilities, many don't. If 60% of queries could be handled locally, that's 60% of inference costs eliminated and 60% of latency removed.

What's your data residency situation? If compliance requirements currently block AI features, edge inference may provide a compliant path. No data leaves the device means no data sovereignty concerns.

What's your connectivity profile? If users are frequently offline, occasionally connected, or in high-latency environments, edge models provide capabilities that cloud-dependent solutions can't match.

For Tech Founders Considering AI Products

Edge-first AI enables business models that cloud-first AI doesn't. Consider:

One-time purchase applications: If your AI runs on-device, you can sell software rather than subscriptions. Users pay once, you have no ongoing inference costs, and the margin structure is entirely different.

Offline-first markets: Emerging markets, industrial applications, and anywhere



Liquid AI Releases LFM2.5-1.2B-Thinking: 1.2B Parameter Reasoning Model Runs Fully Offline in 900 MB on Phones

with unreliable connectivity become viable. These are markets that cloud AI companies structurally can't serve well.

Privacy-as-feature positioning: “Your data never leaves your device” is a marketing message that resonates with increasingly privacy-aware consumers. Edge inference makes that promise credible.

Architectures to Consider

The emerging pattern is **tiered inference**: edge models handle routine requests, cloud models handle complex ones. This requires careful threshold design—you need heuristics or classifiers that route queries appropriately.

A simple approach: run the edge model, measure confidence in its response, and escalate to cloud only when confidence is low. This captures the cost savings of edge inference while preserving cloud capabilities for difficult cases.

More sophisticated approaches use the edge model's response as context for the cloud model, enabling faster convergence on complex queries.

Forward Look: Where This Leads

6 Months: Ecosystem Development

Expect rapid tooling development around models in this size class. The limiting factor for edge AI adoption has been developer experience—getting models running on diverse hardware is currently harder than it should be.

llama.cpp and MLX provide foundations, but purpose-built frameworks for edge reasoning will emerge. Apple's Core ML integration will likely deepen. Android's NNAPI will see more optimization for this model class.

Fine-tuning pipelines for sub-2B models will mature. Today, fine-tuning requires expertise; in six months, services will offer drag-and-drop customization for edge models.

12 Months: Hardware Response

Chip manufacturers will respond to demonstrated demand. Qualcomm, Apple, and



Liquid AI Releases LFM2.5-1.2B-Thinking: 1.2B Parameter Reasoning Model Runs Fully Offline in 900 MB on Phones

AMD will release NPU generations optimized for architectures like LFM2.5. The hybrid convolution-attention design may influence silicon design decisions.

Phone manufacturers will market NPU performance more aggressively once compelling use cases exist. “AI performance” claims will get more specific and more meaningful.

The Larger Trend

This release accelerates a structural shift from centralized to distributed AI inference. The same trend happened with computing generally—from mainframes to PCs to mobile—and AI appears to be following a similar trajectory, just compressed into years rather than decades.

The ultimate destination isn’t either/or. It’s a spectrum with different capability levels deployed at different points in the infrastructure. Foundation models in the cloud, specialized models at the edge, and continuous improvement through federated learning that preserves privacy while enabling collective intelligence.

LFM2.5-1.2B-Thinking is one step on that path—a proof point that edge inference can be genuinely capable, not just a compromise.

What This Means for You

The organizations that benefit most from this release are those that move quickly to test and deploy. Edge AI advantages compound: lower costs fund more features, better latency drives user engagement, and privacy compliance opens new markets.

The organizations that hesitate will find themselves competing against products with structural cost advantages. When your competitor’s inference costs are zero and yours are measured in API calls, their pricing flexibility exceeds yours indefinitely.

This isn’t about replacing cloud AI—it’s about right-sizing deployment to actual requirements. Many applications have been over-provisioned for cloud inference simply because capable edge alternatives didn’t exist. Now they do.

The question isn’t whether edge AI is ready for production—it’s whether



Liquid AI Releases LFM2.5-1.2B-Thinking: 1.2B Parameter Reasoning Model Runs Fully Offline in 900 MB on Phones

your architecture is ready for edge AI.