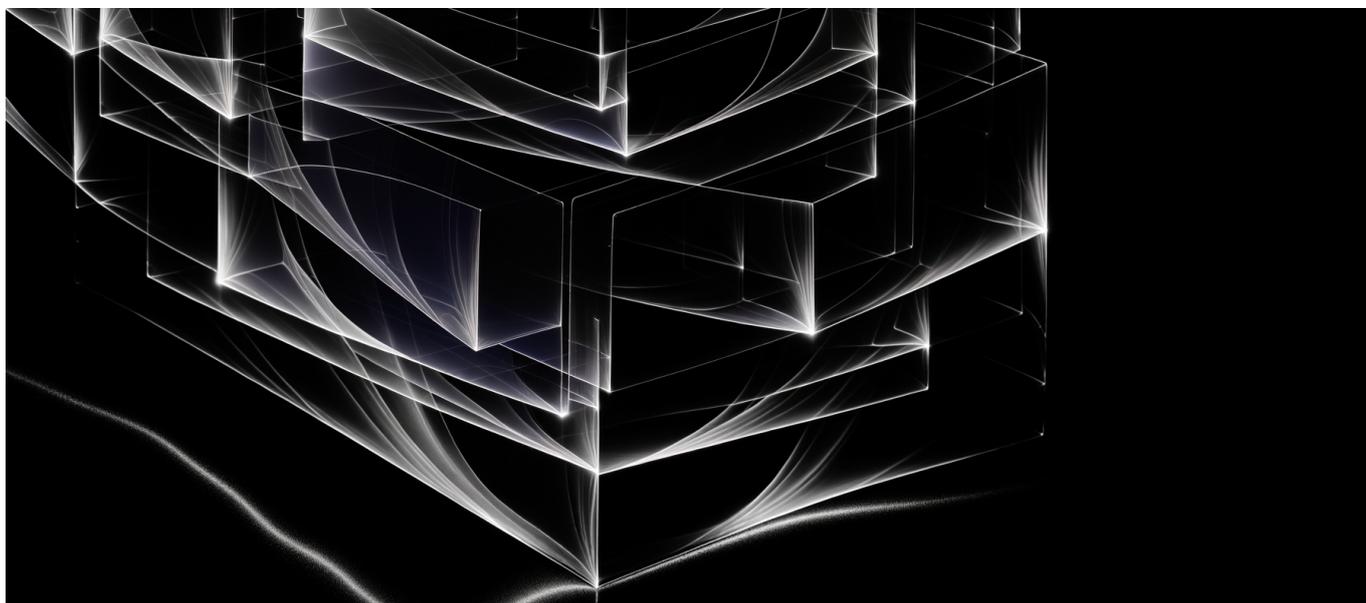




Moonshot AI's Attention Residuals (AttnRes) Achieves Same Training Loss with 1.25× Less Compute—New Transformer Architecture Replaces Fixed Residuals with Softmax Attention



Moonshot AI's Attention Residuals (AttnRes) Achieves Same Training Loss with 1.25× Less Compute—New Transformer Architecture Replaces Fixed Residuals with Softmax Attention

Residual connections have been sacred architecture for nine years—and Moonshot AI just proved they were holding us back. Their new AttnRes method cuts training compute by 25% while matching loss, and the implications ripple through every transformer deployment on the planet.

The News: A Fundamental Assumption Falls

On March 16, 2026, [Moonshot AI announced Attention Residuals \(AttnRes\)](#), a transformer architecture modification that replaces fixed residual connections with



Moonshot AI's Attention Residuals (AttnRes) Achieves Same Training Loss with 1.25× Less Compute—New Transformer Architecture Replaces Fixed Residuals with Softmax Attention

softmax attention between layers. The results are striking: identical training loss with 1.25× less compute, plus a reduction in memory complexity from $O(Ld)$ to $O(Nd)$, where L is layer count, d is hidden dimension, and N is sequence length.

This isn't an incremental optimization. It's a fundamental rethinking of how information flows through transformer networks.

Since Vaswani et al.'s "Attention Is All You Need" paper in 2017, residual connections have been treated as load-bearing walls in transformer architecture. Every major model—GPT-4, Claude, Gemini, Llama—uses the same basic pattern: add the layer's output directly to its input, creating a skip connection that helps gradients flow during training. Nobody questioned it because it worked.

Moonshot AI questioned it anyway. Their insight: why use a fixed, uniform connection when you can let the model learn which earlier layers matter most for each computation?

How AttnRes Actually Works

Traditional residual connections are simple addition. Layer N receives input X , computes some transformation $F(X)$, and outputs $X + F(X)$. The " $X +$ " part is the residual—it's fixed, weighted equally, and applies the same way regardless of what the model has learned.

AttnRes replaces this fixed addition with a learned attention mechanism over all previous layer outputs. Instead of automatically adding the immediate input, each layer attends to a compressed representation of every earlier layer's output using standard softmax attention. The model learns, during training, which historical layer outputs are most relevant for the current computation.

Think of it as giving each layer a selective memory of everything that came before, rather than just a direct line to its immediate predecessor.

The memory complexity reduction from $O(Ld)$ to $O(Nd)$ deserves unpacking. In standard transformers, you need to store activations for each layer (L) times the hidden dimension (d) during the forward pass to enable backpropagation. AttnRes restructures this so memory scales with sequence length (N) times hidden dimension (d), decoupling memory requirements from model depth. For a 100-layer model, this distinction becomes substantial.



Moonshot AI's Attention Residuals (AttnRes) Achieves Same Training Loss with 1.25× Less Compute—New Transformer Architecture Replaces Fixed Residuals with Softmax Attention

The Computational Trade-off

You might wonder: doesn't adding attention between layers increase compute? It does add some overhead to each forward pass. But AttnRes achieves its 1.25× compute reduction through faster convergence—the model reaches the same loss in fewer training steps because information flows more efficiently.

This is counterintuitive. Adding computation per step while reducing total computation means each step is doing more useful work. The learned attention over previous layers apparently helps the model coordinate its internal representations more effectively than fixed skip connections.

Why This Matters More Than Most Architecture Papers

The timing amplifies the significance. Training compute costs are the binding constraint on AI development right now. AMI Labs [raised \\$1.03 billion on March 10, 2026](#) for world model research—and a significant chunk of that will go to GPU clusters. Any genuine reduction in training compute translates directly to either cost savings or capability gains.

A 25% compute reduction doesn't sound dramatic until you multiply it by the scale involved. Training a frontier model costs \$50-100 million in compute alone. A 1.25× efficiency gain means you either save \$10-20 million per training run, or you train a model 25% longer (potentially larger or more capable) for the same budget.

Compounding this, the Cerebras-AWS partnership announced on the same day, [March 16, deployed CS-3 systems achieving 5× token throughput](#). AttnRes-style architecture improvements combined with specialized hardware improvements could deliver multiplicative gains. A 1.25× from architecture times a 5× from hardware yields over 6× total improvement—the kind of shift that changes what's economically feasible.

The Memory Advantage Is Underrated

Most coverage has focused on the compute reduction, but the memory complexity change from $O(Ld)$ to $O(Nd)$ deserves equal attention.



Moonshot AI's Attention Residuals (AttnRes) Achieves Same Training Loss with 1.25× Less Compute—New Transformer Architecture Replaces Fixed Residuals with Softmax Attention

Deep transformers have always faced a memory wall. Training a 200-layer model requires storing intermediate activations for all 200 layers during the forward pass to enable gradient computation during backpropagation. This forces practitioners to either use gradient checkpointing (trading compute for memory) or limit model depth.

AttnRes breaks this coupling. Memory now scales with sequence length rather than depth, meaning you can train much deeper models on the same hardware. Or equivalently, train the same depth model on cheaper hardware with less memory.

For organizations running on consumer or mid-tier GPUs rather than H100 clusters, this memory efficiency might matter more than the compute savings. It expands who can meaningfully participate in model development.

Technical Deep Dive: Architecture and Training Dynamics

Let's get specific about what AttnRes changes in the forward pass.

In a standard transformer block, you have:

$$\text{Output} = \text{LayerNorm}(X + \text{Attention}(X) + \text{FFN}(X))$$

The "X +" terms are residual connections, adding the input directly to the output of each sub-layer.

AttnRes modifies this to:

$$\text{Output} = \text{LayerNorm}(\text{AttnRes}(H_1 \dots H_n) + \text{Attention}(X) + \text{FFN}(X))$$

Where $H_1 \dots H_n$ are the outputs of all previous layers, and AttnRes is a lightweight attention mechanism that learns to weight them.

The AttnRes attention uses the current layer's representation as the query, and compressed representations of all previous layer outputs as keys and values. This compression is essential—without it, the memory overhead would be prohibitive. The compression appears to use a projection down to a lower-dimensional space, similar to how multi-head attention reduces dimensionality within each head.



Moonshot AI's Attention Residuals (AttnRes) Achieves Same Training Loss with 1.25× Less Compute—New Transformer Architecture Replaces Fixed Residuals with Softmax Attention

Why This Helps Gradient Flow

Residual connections were originally introduced to solve the vanishing gradient problem in very deep networks. By providing a direct path for gradients to flow backward, they made training 100+ layer networks feasible.

AttnRes potentially improves on this by making the gradient paths adaptive rather than fixed. When the model learns that layer 47's output is particularly relevant for layer 98's computation, gradients can flow directly along that learned connection. The fixed residual pattern forces gradients through every intermediate layer even when many of those layers aren't contributing meaningfully to a particular computation.

This selective gradient flow could explain the faster convergence. The model wastes less capacity routing gradients through irrelevant pathways.

Comparison to Prior Work

This isn't the first attempt to improve on residual connections. DenseNet (2017) connected every layer to every other layer with concatenation rather than addition. Highway Networks (2015) added learned gating to residual connections.

AttnRes differs in using attention rather than concatenation or gating. Attention provides a more flexible learned weighting that can dynamically adjust based on the input, rather than using fixed gates. It's also more parameter-efficient than DenseNet's approach, which required each layer to handle concatenated inputs of growing dimension.

The Ai2 Olmo Hybrid model released in early March 2026 showed [2× data efficiency with 49% fewer tokens on MMLU](#) through a different architectural modification (hybridizing attention with state-space models). The timing suggests we're in a period of rapid architectural exploration, with multiple teams finding significant efficiency gains by questioning long-standing assumptions.

What Most Analysis Gets Wrong

The obvious narrative is "Moonshot AI found free efficiency gains, everyone should adopt AttnRes immediately." This misses several important nuances.



Moonshot AI's Attention Residuals (AttnRes) Achieves Same Training Loss with 1.25× Less Compute—New Transformer Architecture Replaces Fixed Residuals with Softmax Attention

The Overhyped Part: Universal Applicability

A 1.25× compute reduction on Moonshot AI's specific benchmark configuration doesn't automatically transfer to every training setup. The benefit depends on:

Model depth: The memory advantage scales with layer count. For shallow models (under 24 layers), the overhead of the AttnRes mechanism might not pay for itself.

Sequence length: The new $O(Nd)$ memory complexity means very long sequences could actually increase memory requirements compared to shallow models with the old $O(Ld)$ formula. The crossover point depends on the specific L , N , and d values.

Training configuration: Batch sizes, learning rates, and optimizer choices all interact with architectural changes. AttnRes might require hyperparameter re-tuning to achieve advertised gains.

Task type: The benefits measured on language modeling might not transfer directly to vision, multimodal, or other transformer applications.

The Underhyped Part: Inference Implications

Almost no coverage has addressed what AttnRes means for inference. The attention over previous layers adds computation to every forward pass. Training efficiency gains come from faster convergence, but inference doesn't involve convergence—you're just running forward passes.

If AttnRes adds even 10% overhead per forward pass, that matters enormously at inference scale. A model serving millions of queries multiplies that overhead across every request. The training efficiency gains are a one-time benefit; inference overhead is a recurring cost.

This suggests AttnRes might be most valuable for training large models that then get distilled, quantized, or otherwise optimized for inference. You'd capture the training efficiency during development, then eliminate the inference overhead through model compression.

Alternatively, there might be inference-specific optimizations possible—caching certain layer attention patterns, approximating the AttnRes attention with learned fixed weights after training, or other techniques not yet explored.



Moonshot AI's Attention Residuals (AttnRes) Achieves Same Training Loss with 1.25× Less Compute—New Transformer Architecture Replaces Fixed Residuals with Softmax Attention

The Actually Interesting Question: What Does This Imply About Current Models?

If AttnRes achieves identical loss with less compute, what does that say about existing models trained with fixed residuals?

One interpretation: current models are slightly overtrained. They've used more compute than strictly necessary to reach their capability level because the architecture wasn't optimal.

Another interpretation: current models have learned suboptimal internal representations because the fixed residual connections constrained how information could flow. An AttnRes model might develop qualitatively different internal structure, even when achieving the same loss.

The second interpretation is more interesting. It suggests that loss-equivalent models might behave differently on out-of-distribution inputs, have different failure modes, or generalize differently. The loss metric captures average performance on training-like data but doesn't fully characterize model behavior.

Practical Implications: What You Should Actually Do

For different reader profiles, the action items differ:

If You're Training Large Models (Frontier Labs, Well-Funded Startups)

Allocate engineering time to replicate AttnRes on a smaller scale. Run an ablation study comparing your current architecture to an AttnRes variant on a 1B-7B parameter model. If you see convergence improvements that match Moonshot's claims, plan integration into your next large training run.

The 25% compute savings compound with scale. At \$50M training budgets, this is worth significant engineering investment to validate.

Watch for Moonshot AI's technical report or paper. The announcement provides results but presumably lacks implementation details. The full publication will answer



Moonshot AI's Attention Residuals (AttnRes) Achieves Same Training Loss with 1.25× Less Compute—New Transformer Architecture Replaces Fixed Residuals with Softmax Attention

questions about hyperparameter sensitivity, failure modes, and edge cases.

If You're Building Applications on Existing Models

No immediate action required. Your API-accessed models won't change based on this announcement, and the providers will handle architecture decisions.

However, understand that the next generation of foundation models might use AttnRes or similar techniques. If you're planning long-term integrations, recognize that inference characteristics could shift—potentially different latency profiles, different behavior on edge cases, different failure modes.

If You're Working on Model Efficiency

AttnRes joins a growing family of architectural efficiency improvements from early 2026. Combine this with:

- Mixture of Experts (MoE) routing improvements
- Hybrid attention/state-space architectures (like Ai2's Olmo Hybrid)
- Specialized inference kernels for new hardware (like Cerebras CS-3)

The pattern suggests that significant efficiency gains remain available through architectural innovation. The transformer architecture isn't optimized; it's merely adequate. Teams that systematically explore the design space are finding 25-100% improvements that were available all along but assumed away.

If You're Making Vendor/Hardware Decisions

The memory complexity reduction ($O(Ld)$ to $O(Nd)$) changes the optimal hardware configuration for large model training. Models using AttnRes need less memory per layer but might need more compute per forward pass.

This favors hardware with high compute-to-memory ratios, assuming AttnRes becomes standard. It also slightly reduces the advantage of extremely high-memory configurations, since the memory wall has moved.

Don't make immediate purchasing decisions based on this, but factor it into 12-18 month planning horizons. If AttnRes or similar approaches become standard, the optimal GPU/accelerator balance shifts.



Moonshot AI's Attention Residuals (AttnRes) Achieves Same Training Loss with 1.25× Less Compute—New Transformer Architecture Replaces Fixed Residuals with Softmax Attention

The Competitive Landscape Shifts

Moonshot AI is a China-based company, and this announcement has geopolitical undertones worth acknowledging.

The US has implemented export controls on advanced AI chips to China. These controls aim to slow Chinese AI development by restricting access to hardware. Architectural efficiency improvements like AttnRes partially counter this strategy—if you can do more with less compute, hardware restrictions have less impact.

A 25% efficiency gain doesn't fully offset restricted access to H100s, but it helps. More importantly, it suggests that Chinese AI labs are investing heavily in algorithmic efficiency research, potentially as a strategic response to hardware constraints.

For US/European organizations, this creates a complicated dynamic. AttnRes appears to be a genuine technical improvement that would benefit anyone who adopts it. But adopting it means building on Chinese AI research, with potential IP complications and security considerations for some applications.

The pragmatic approach: evaluate AttnRes purely on technical merit, but maintain awareness of the broader context. If the technique is valid, it will be independently verified and implemented by multiple labs. Wait for independent reproduction before major production commitments.

Where This Goes: A 6-12 Month View

Near Term (3-6 Months)

Expect rapid replication attempts from major labs. If AttnRes delivers as claimed, you'll see papers from Google, Meta, Anthropic, and others either confirming the results or explaining why they don't replicate in specific configurations.

Moonshot AI's credibility will be established or damaged based on independent reproduction. Successful replication likely leads to fast adoption; failed replication relegates this to the "interesting but not general" category.

Open-source implementations will emerge. The technique isn't obviously patent-blocking, and the research community moves fast. HuggingFace implementations



Moonshot AI's Attention Residuals (AttnRes) Achieves Same Training Loss with 1.25× Less Compute—New Transformer Architecture Replaces Fixed Residuals with Softmax Attention

within 60-90 days are plausible if the technical report provides sufficient detail.

Medium Term (6-12 Months)

If AttnRes proves out, expect the next generation of frontier models (GPT-5 class, Claude 4 class) to incorporate it or similar learned-residual approaches. The efficiency gain is too large to leave on the table at frontier training scales.

Derivative research will explore variations: different compression schemes for the layer history, different attention mechanisms, combination with other efficiency techniques.

The more interesting possibility: AttnRes proves to be one instance of a broader principle. Perhaps there are other fixed architectural choices in transformers—layer normalization placement, attention head patterns, FFN expansion ratios—that could similarly be replaced with learned alternatives for additional efficiency gains.

The transformer architecture was designed in 2017 for the hardware and scale available then. We're training models 1000× larger on hardware several generations ahead. A systematic re-examination of every architectural choice against modern constraints seems overdue, and AttnRes suggests such re-examination yields significant returns.

The Speculative Upside

If multiple 25%-scale efficiency improvements can be combined—better residuals, hybrid architectures, improved attention patterns, optimized FFN structures—the compound effect could exceed 2× total efficiency gain.

A 2× efficiency gain means current hardware can train models twice as capable, or current capability levels become accessible to organizations with half the budget.

This acceleration isn't about faster progress toward AGI or other dramatic claims. It's about democratization and cost reduction. The difference between "\$50M to train a frontier model" and "\$25M to train a frontier model" is the difference between a handful of organizations being able to compete and dozens being able to compete.

More competition generally produces faster iteration and better outcomes. AttnRes,



Moonshot AI's Attention Residuals (AttnRes) Achieves Same Training Loss with 1.25× Less Compute—New Transformer Architecture Replaces Fixed Residuals with Softmax Attention

if validated and combined with other efficiency improvements, contributes to this dynamic.

What This Says About the Field

The residual connection has been a fixed assumption for nine years. Millions of engineering hours have been spent optimizing transformers while treating residuals as immutable. And a relatively straightforward modification—using attention instead of addition—yields 25% efficiency gains.

This should make everyone slightly uncomfortable. What other “load-bearing walls” in neural network architecture are actually optional? What other assumed-necessary components are just locally optimal defaults that nobody questioned?

The lesson isn't specific to AttnRes. It's that architectural exploration remains under-invested relative to the returns available. The transformer isn't optimized; it's the first thing that worked at scale, and we've been iterating within its constraints rather than questioning them.

Moonshot AI's contribution isn't just a specific technique. It's evidence that the design space remains large and underexplored.

The teams that systematically question inherited assumptions—rather than optimizing within them—will define the next generation of AI architecture, and Moonshot AI just demonstrated the returns available to those who ask whether the residual connection was really the best we could do.