

# Why Hybrid Architectures Matter

Combining Transformers and State-Space Models for  
Next-Generation Language Modeling

Bryan Leonard & Brandyn Leonard

*Qira LLC, Maricopa, AZ*

---

## Abstract

*Modern language models overwhelmingly rely on a single architectural paradigm: the Transformer decoder. While attention mechanisms excel at capturing local syntax and retrievable factual associations, they struggle with slow-changing discourse structure, topic continuity, and long-range sequential dependencies. State-space models (SSMs) offer complementary strengths—linear-time recurrence, continuous state evolution, and natural handling of sequential dynamics—but lack the precise token-to-token retrieval that attention provides. We argue that neither architecture alone is sufficient for next-generation language models, and present evidence from the LOLM (Latent Order Language Model) architecture that a principled hybrid approach, combining a Transformer surface decoder with a Mamba-style latent SSM through a learned per-dimension manifestation gate, produces emergent properties that neither component exhibits in isolation. We report a dependency inversion phenomenon at 1.57B parameters where the minority latent path (30% gate weight) becomes more structurally essential than the majority surface path, suggesting that hybrid architectures enable qualitatively different internal representations.*

## 1. Introduction

The Transformer architecture has dominated language modeling since its introduction in 2017. Scaled to hundreds of billions of parameters and trained on trillions of tokens, Transformer-based models have achieved remarkable performance across diverse language tasks. Yet the architecture has well-known limitations: quadratic attention complexity with sequence length, difficulty with very long-range dependencies, and no explicit mechanism for tracking slow-changing discourse variables like topic, intent, or rhetorical mode.

State-space models, particularly the selective SSM architecture introduced by Mamba, offer an alternative paradigm. With linear-time complexity, content-aware state filtering, and continuous recurrent dynamics, SSMs naturally capture sequential structure that attention handles inefficiently. Recent work has shown that Mamba-3B matches the perplexity of a Transformer-6B while being significantly cheaper to run at inference time.

However, pure SSMs struggle with tasks requiring precise memory retrieval and in-context learning—capabilities where attention excels. This complementary failure pattern suggests that the optimal architecture is neither pure Transformer nor pure SSM, but a hybrid that combines both.

## 2. The Attention-Recurrence Tradeoff

Attention and recurrence represent fundamentally different inductive biases for sequence modeling. Attention computes pairwise relationships between all positions, enabling direct information routing from any token to any other. This makes attention powerful for factual retrieval ("what was mentioned earlier?") and syntactic agreement across distance. However, attention is position-agnostic without explicit positional encoding, treats each layer independently, and has no intrinsic notion of sequential state evolution.

Recurrence, by contrast, maintains a continuous hidden state that evolves through time. Selective SSMs like Mamba make this evolution input-dependent: the discretization step, input projection, and output projection all adapt to current context. This creates a content-aware filter that naturally tracks slow-changing variables—precisely the discourse-level structure that attention handles poorly.

The key insight motivating hybrid architectures is that these two capabilities are not redundant. A Transformer excels at "what" questions (what token should come next, given precise retrieval of relevant context). An SSM excels at "where" questions (where in the discourse are we, what mode is the text operating in, what sequential patterns are unfolding). A hybrid model can delegate each question to the component best equipped to answer it.

### 3. Existing Hybrid Approaches

Several recent architectures combine Transformers and SSMs. Jamba interleaves Mamba and Transformer layers with a mixture-of-experts architecture. Griffin alternates between local attention and recurrent layers. Samba places Mamba, sliding-window attention, and MLP layers in a specific repeating pattern. These approaches share a common design philosophy: fixed interleaving of heterogeneous layers, where the mixing ratio is an architectural hyperparameter rather than a learned quantity.

LOLM takes a different approach. Rather than interleaving layers, it runs a complete Transformer decoder and a complete SSM core in parallel, producing separate hidden representations that are combined through a learned per-dimension gate. This design has three advantages: (1) each pathway operates with its full depth and capacity, (2) the mixing ratio is learned per-dimension and per-position rather than fixed, and (3) the gate provides an interpretable window into how the model allocates computation between surface and latent representations.

### 4. Per-Dimension Gating Changes the Equation

Most hybrid architectures use scalar mixing or fixed ratios to combine their components. LOLM's manifestation gate is a 2-layer MLP that produces a gating vector  $g$  in  $[0,1]^d$ , where  $d$  is the model dimension. This means each feature dimension independently decides how much information to draw from the surface decoder versus the latent SSM.

The gate takes as input all four information streams: the surface hidden state  $h$ , the latent state  $z$ , the memory readout  $m$ , and the regime embedding  $r$ . This fully-informed arbitration means the gating decision is contextual—the model can shift toward latent representations when processing discourse-level transitions and toward surface representations for local syntactic completion.

At convergence on the 304M model, the gate settles at approximately 0.83, meaning roughly 17% of the fused representation comes from the latent SSM. At 1.57B, the gate settles lower at approximately 0.71,

allocating 29% to latent. The model learns to use more latent capacity at larger scale—a finding consistent with the hypothesis that latent structure becomes more valuable as models grow.

## 5. Dependency Inversion: The Core Finding

The most striking result from LOLM is what we term dependency inversion. At 1.57B parameters, forcing the gate to 1.0 (surface only, removing the 30% latent contribution) causes perplexity to explode to 485 million—a complete model collapse. Forcing the gate to 0.0 (latent only, removing the 70% surface contribution) produces perplexity of 56,130—bad, but the model still functions as a language model.

The asymmetry is extreme: the surface-only collapse is 8,645 times worse than the latent-only collapse. The minority contributor is more structurally essential than the majority contributor. This is not how auxiliary modules behave in standard architectures. When you add a small auxiliary component to a Transformer and then remove it, the base model degrades gracefully. Here, the base model has reorganized itself to depend on the auxiliary path so completely that it cannot function without it.

This finding suggests that hybrid architectures, when trained with appropriate gradient isolation and per-dimension gating, enable a qualitatively different mode of internal organization—one where the model distributes essential computation across heterogeneous pathways rather than treating one as primary and the other as supplementary.

## 6. Implications for Architecture Design

If dependency inversion is a general property of well-trained hybrid architectures (rather than an artifact of LOLM's specific design), it has significant implications for how we think about scaling language models. It suggests that adding heterogeneous computational pathways—SSMs, memory systems, discrete codebooks—is not merely about adding capacity. It is about enabling the model to discover internal structures that no single pathway can represent alone.

The practical question is whether this translates to improved downstream task performance, not just lower perplexity. Our controlled comparison at 1.57B shows approximately 2x compute efficiency and 25.9% lower perplexity at matched training budget. Downstream benchmark evaluations (MMLU, HellaSwag, ARC) are in progress and will be reported in future work.

## 7. Conclusion

Neither Transformers nor state-space models alone capture the full structure of natural language. Transformers excel at local precision and retrieval; SSMs excel at sequential dynamics and continuous state tracking. Hybrid architectures that combine both through learned, contextual gating can achieve emergent integration properties—including dependency inversion—that suggest a fundamentally different mode of learned representation. The evidence from LOLM indicates that this is not merely additive benefit but a qualitative shift in how models organize their internal computation. The hybrid paradigm deserves deeper investigation as a path toward more capable and efficient language models.

Qira LLC — Leonard & Leonard, 2026