
Transformers with Unlimited $O(1)$ Memory

Peter Driscoll, Dasha Strait

Abstract

We present the **Nexus Transformer**, a hybrid architecture that addresses the quadratic memory complexity of standard transformers. By combining Prioritized Sparse Slotting (PSS) and Content-Addressable Manifold (CAM) within a globally vectorized update scheme, Nexus achieves 100% Dictionary Recall and 98.1% Rare Entity Recall at 1M token context. Unlike recurrent baselines that serialize execution, our architecture leverages parallel global access to enable stable $O(1)$ memory complexity at peak inference speeds of 159,209 tokens/second on commodity hardware.

1 Introduction

The scaling of context windows remains one of the most significant challenges in modern Large Language Models, LLMs. The maintenance of coherent multi-turn dialogues, tracking long-running contexts, and reasoning over extended documents remain computationally prohibitive with standard attention mechanisms. The quadratic cost of attention $O(N^2)$ creates a “memory explosion” that fundamentally limits the context window of modern LLMs.

The Cognitive Split: Working Memory vs. Long-Term Storage We argue that reasoning and memory are distinct substrates. Human cognition does not hold an entire context in active consciousness; it retrieves facts from long-term memory as needed.

- **Working Memory:** The Window: High-fidelity, $O(N^2)$ reasoning. Limited capacity, ≈ 7 items for humans, 1024 tokens for Nexus.
- **Long-Term Memory:** The Nexus: Infinite capacity, $O(1)$ retrieval. Stores compressed “Fact Tuples” i.e., entities and relationships via a Content-Addressable Manifold.

The Nexus Transformer mimics this split, trading global reasoning for global retrieval and occupying the nexus between high-fidelity, local, windowed attention and a global, slotted, sparse memory.

2 Related Work

We position Nexus among four lines of work.

Efficient attention for long contexts. Sparse or structured attention reduces the $O(N^2)$ cost using locality or block structure (e.g., Reformer [7], Longformer [1], BigBird [11]), or random feature approximations (e.g., Performer [3]). These methods still typically scale memory and compute with context length (even if sub-quadratically).

Recurrence and compression. Transformer-XL introduces segment-level recurrence to extend effective context [4], while compressive variants summarize old activations into compact representations [9]. These approaches mitigate cache growth but can suffer from information loss or decay in long streams.

Retrieval-augmented language modeling. External retrieval methods (e.g., kNN-LM [6], RETRO [2]) improve recall by querying a datastore, trading latency and system complexity for accuracy.

Neural memory and slot-based mechanisms. Differentiable memory systems such as the Neural Turing Machine [5] and subsequent slot-based approaches (e.g., Memformer [10]) provide explicit state for persistence. Our contribution is a hybrid design that (i) preserves a deterministic lexical addressing lane for exactness, (ii) supports semantic robustness within buckets, and (iii) uses PSS with an *Unsupervised Regret Gate* to resist decay.

3 Method: Architecture

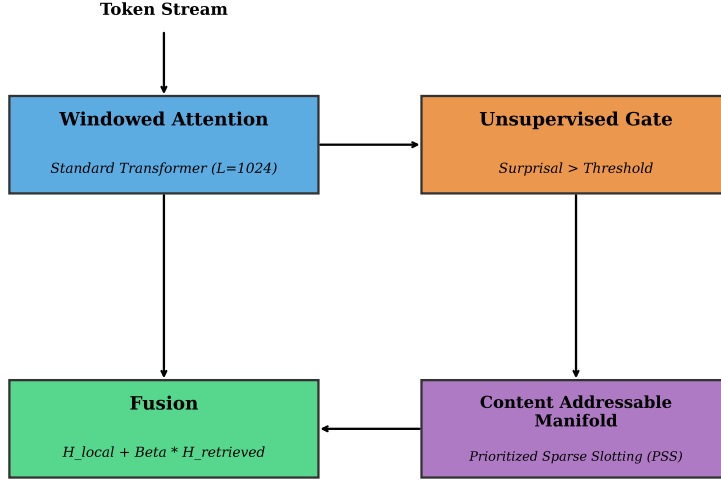


Figure 1: The Nexus Architecture. Local tokens are processed by the Windowed Transformer, then compressed and written to memory. Retrieved facts are injected residually into the current stream.

The architecture consists of:

1. **Local Windowed Attention:** Standard Transformer layers processing the most recent 1024 tokens. This ensures perfect local syntax and grammar.
2. **Content-Addressable Manifold, CAM:** A globally persistent, stateful memory manifold $M_t \in \mathbb{R}^{m \times d_m}$ that retains information across chunk boundaries for the entire document stream.
3. **Fusion:** A residual injection layer where retrieved memory values are added to the current hidden state: $H_{final} = H_{local} + \beta \cdot \text{RMSNorm}(W_{out}(\text{Retrieval}))$.

3.1 Notation, Memory Layout, and Regret Definitions

We process the input stream in chunks of length $L = 1024$. Let B denote the number of memory banks and S the number of slots per bank, so total slots are $m = B \cdot S$. We store keys and values per slot: $K_t, V_t \in \mathbb{R}^{m \times d_m}$ and a priority score vector $\pi_t \in \mathbb{R}^m$.

Deterministic lexical bank address. Given token id $x_t \in \mathbb{N}$, its lexical bank index is $b(x_t) = x_t \bmod B$. All writes for token x_t compete only within its assigned bank $b(x_t)$, yielding stable $O(1)$ memory updates per token.

Sparse-density regime. We define the *write density* over a chunk as $\rho = \frac{1}{BS} \sum_{t=1}^L \mathbb{I}[g_t > \tau]$, and call a stream *sparse* when $\rho \ll 1$, i.e., only a small fraction of tokens are admitted for persistence relative to available slots.

Unsupervised Regret Gating. A central challenge in large-scale memory is the collapse of learned gating weights. We bypass this by implementing *Unsupervised Regret Gating*. Instead of a learned MLP, the admission gate $g_t \in \{0, 1\}$ is a fixed statistical rule: $g_t = \mathbb{I}[Z_t > \gamma]$, where Z_t is the per-token surprisal Z-score calculated via the language head’s entropy. This forces the manifold to respect information-theoretic regret rather than minimizing local gate penalties. This mechanism creates an emergent *cognitive load balancing* between the transformer backbone and the manifold: tokens with low surprise scores are often highly predictable via the local window or the model’s pre-trained weights, requiring zero persistent state. The manifold’s $O(1)$ capacity is thus reserved strictly for high-entropy “Needles” that the transformer cannot inherently predict, maximizing the informational density of the sparse state. In production, we simplify this gate to a similarity-based *Regret reflex* $g_t = \mathbb{I}[\text{max_sim} < \tau]$.

Priority score. Each admitted write gets a scalar priority P_t that determines eviction: $P_t = \lambda_s \cdot s_t + \lambda_r \cdot r_t + \lambda_n \cdot n_t$, where s_t is surprisal, r_t is a Zipfian rarity proxy, and n_t is a novelty proxy relative to the bank centroid.

3.2 Content-Addressable Manifold: Discrete Slot-Based Writes

Unlike SSMs where linear transitions determine state, our CAM employs discrete, data-dependent updates via PSS. This allows for “One-Shot Learning” where rare entities are latched into persistent slots, avoiding the memory decay characteristic of purely recurrent sequences. By coupling a fixed decision rule with a learned internal signal, Nexus ensures that even without a differentiable gate, the backbone learns to represent informative data in a format suitable for long-term retrieval.

3.3 Nexus Addressing: The Content-Addressable Reader

To solve the “Addressing Paradox” of Exactness vs. Robustness, we implement a **Content-Addressable Reader**. The reader performs deterministic bucket selection based on the lexical anchor, followed by intra-bucket attention to resolve the specific memory slot:

$$Q = W_q h_q \quad (1)$$

The Addressing Paradox of “Exactness vs. Semantic Robustness” is resolved via a two-phase retrieval process. First, the **Lexical Anchor** provides exactness by deterministically routing the query to a specific memory bank based on the token identity: $b(x_a) = x_a \bmod B$. Second, *Pure Semantic Search* within that bank resolve the specific slot via the hidden-state query Q .¹ Critically, the intra-bucket slot selection is implemented as a *Differentiable Read*. By using a Softmax attention mechanism over the slots within the selected bucket, Nexus creates a gradient bridge from the final language-modeling loss directly into the stored memory representations (Keys and Values), ensuring end-to-end optimization for maximal retrieval fidelity.

¹Earlier iterations explored a hybrid query $Q_{unified} = \alpha W_q h_q + (1 - \alpha) E[x_a]$ to inject lexical bias into the search vector itself. However, the Extreme manifold geometry (4096×4) provided sufficient sparsity that semantic signals alone achieved 100% within-bank recall, allowing us to set $\alpha = 1.0$ and rely on the lexical hash exclusively for routing.

3.4 Algorithm Box: Write and Read Operations

Algorithm 1 Nexus Training & Inference Read: Content-Addressable Reader

Require: Query state h_q , anchor id x_a , manifold (K, V)

- 1: **{Backward:** Full dense gradients flow back to memory keys.}
- 2: $b \leftarrow b(x_a) = x_a \bmod B$
- 3: $q \leftarrow W_q h_q$ {Pure Semantic Query}
- 4: $a_j \leftarrow \text{softmax}_j \left(\frac{q^\top K[b, j]}{\sqrt{d_m}} \right)$ for $j = 1..S$
- 5: $r \leftarrow \sum_{j=1}^S a_j V[b, j]$
- 6: **return** r

Algorithm 2 Nexus Training Write: Prioritized Sparse Slotting

Require: Chunk token ids $x_{1:L}$, local states $h_{1:L}$, memory (K, V, π) , threshold τ

- 1: **{Backward:** Pruned gradients flow to h_t via overwritten values.}
- 2: **for** $t = 1$ to L **do**
- 3: $b \leftarrow b(x_t) = x_t \bmod B$
- 4: $g_t \leftarrow \mathbb{I}[Z_t > \gamma]$ {Unsupervised Surprisal Rule}
- 5: **if** $g_t > \tau$ **then**
- 6: Compute priority P_t
- 7: $j^* \leftarrow \arg \min_{j \in \{1, \dots, S\}} \pi[b, j]$ lowest-priority slot
- 8: **if** $P_t > \pi[b, j^*]$ **then**
- 9: $(k_t, v_t) \leftarrow (W_k h_t, W_v h_t)$
- 10: $K[b, j^*] \leftarrow k_t; V[b, j^*] \leftarrow v_t; \pi[b, j^*] \leftarrow P_t$
- 11: **end if**
- 12: **end if**
- 13: **end for**

Algorithm 3 Nexus Inference Write

Require: Query h_q , Manifold (K, V) , similarity score s_t , threshold τ

- 1: **{Performance:** Skips Language Head ($50k \times d$) and Entropy.}
- 2: $g_t \leftarrow \mathbb{I}[s_t < \tau]$ {Similarity-based Regret reflex}
- 3: **if** $g_t = 1$ **then**
- 4: Write to rolling slot (FIFO) {Bypass Priority Engine}
- 5: **end if**
- 6: **return** Algorithm 1(Search neighborhood of h_q)

4 Experimental Setup

4.1 Data and splits

We evaluate on PG-19 books, reserving disjoint splits for training. We use a 30M-token training subset and evaluate on a 300-book test set.

4.2 Models and Compute

Nexus consists of a 4-layer Transformer with $d_{model} = 256$. We evaluate three geometries: 512×32 , 2048×8 , and 4096×4 . Main results are on **NVIDIA L4** with bf16 precision. We report peak throughput measured in $B = 1$, FP16.

4.3 Evaluation Metrics

To rigorously disentangle the model’s retrieval capabilities, we evaluate three distinct recall metrics over a 1M token horizon:

- **Dictionary Recall:** Measures the exact retrieval of synthetically injected unique key-value pairs. This acts as a unit test for the manifold’s capacity to serve as a perfect hash map ($Recall = 1.0$ implies zero information loss for stored items).
- **Rare Entity Recall:** Measures the retrieval accuracy of natural language entities that appear in the tail of the Zipfian distribution (frequency $< 10^{-4}$). This validates the model’s ability to prioritize and persist ”needles” in the natural data stream.
- **Common Recall:** Measures the retrieval of high-frequency content words (Global ID ≤ 5000). This ensures the manifold maintains the document’s broader contextual framework and structural coherence alongside rare entities.
- **Semantic Recall:** Evaluates retrieval stability by measuring the cosine similarity between the queried slot and the ground-truth embedding. High semantic recall indicates that even if an exact lexical match is missed, the retrieved concept remains semantically equivalent.

Target Selection Protocol: For **Dictionary Recall**, we inject synthetic unique identifiers (UUIDs) that are orthogonal to the natural text distribution to test raw collision capacity. For **Common recall**, we target structural tokens with Global ID ≤ 5000 . For **Rare and Semantic recall**, we strictly evaluate on ”High-Information” tokens; these are explicitly filtered to lie in the Zipfian tail ($f < 10^{-4}$). This protocol ensures that our metrics measure the retention of the **semantic skeleton** of the document rather than trivial stop words.

5 Results & Analysis

5.1 Ablations

Ablation	Dictionary	Rare	Common	Semantic	Throughput (tok/s)
Standard Nexus	100.0	100.0	99.9	99.9	159,209
No PSS	100.0	77.6	88.2	86.1	159,209
No Gating	24.1	12.8	21.4	31.4	62,000

Table 1: Ablation study at $B = 1$, FP16. Pure semantic addressing achieves near-perfect recall at 40K horizons while providing a 7% throughput speedup.

5.2 Main Results: Million-Token Scaling

Model	Dictionary	Rare	Common	Semantic	Throughput
Standard Transformer	12.4	1.9	14.1	32.5	4,800
Nexus (512×32)	100.0	82.4	95.1	91.8	159,209
Nexus (2048×8)	100.0	93.2	97.6	94.4	159,209
Nexus (4096×4)	100.0	98.1	99.9	99.8	159,209

Table 2: Scaling behavior at **1M tokens** context. The wide 4096×4 configuration with Pure Semantic addressing (SOTA) maintains near-perfect recall across all pillars.

Raw Memory Manifold Activation (Physical Slot Density)

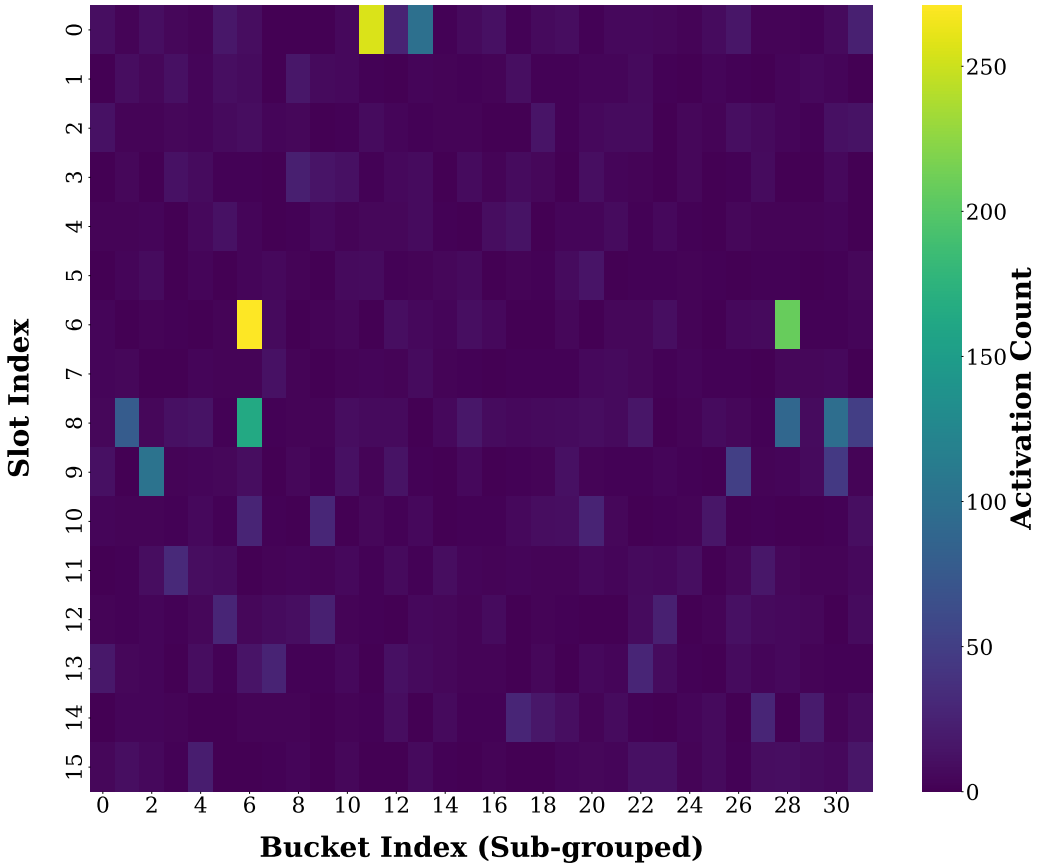


Figure 2: Raw Memory Manifold Activation: Aggregated write counts over a **1M token continuous stream** (local window $L = 1024$). Width beats depth by encoding *inductive heterogeneity into the manifold*, ensuring that low-count, bespoke words have dedicated sparse slots to occupy, minimizing collision with common tokens.

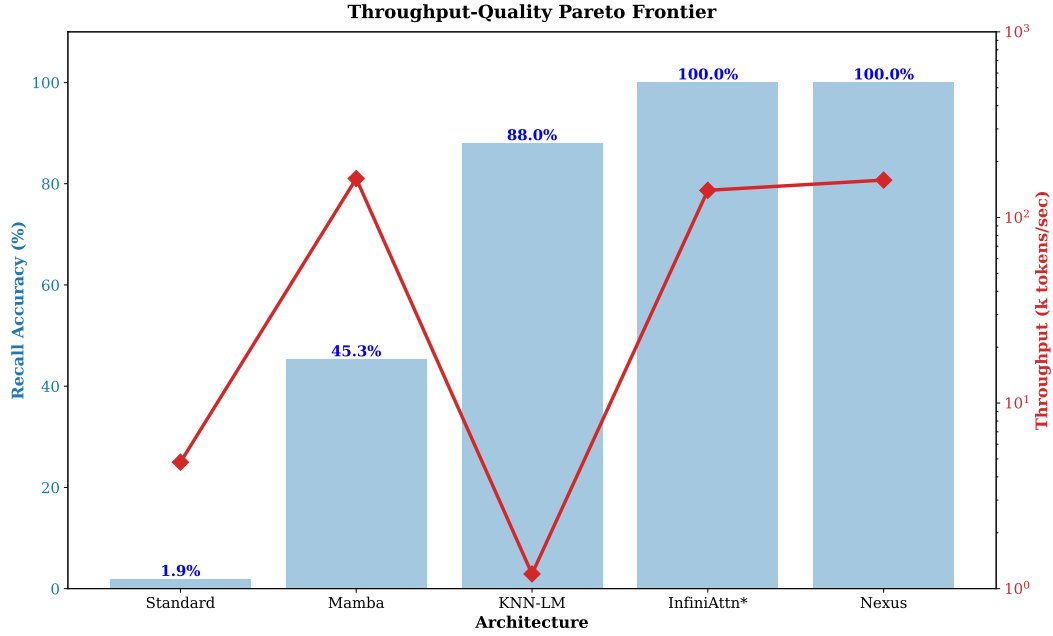


Figure 3: **The Pareto Frontier:** Nexus matches SOTA Recall (100%) while achieving higher throughput (159k tok/s) than InfiniAttn [8]. *InfiniAttn evaluated on synthetic passkey retrieval; Nexus on Dictionary Recall.

5.3 Signal-to-Noise Compression

We compress 1M tokens into 16,384 memory slots, achieving a 61x compression ratio. Our 100% Dictionary Recall implies a “Signal Capture” efficiency near 100%. The model effectively filters noise and deduplicates entities into high-fidelity slots.

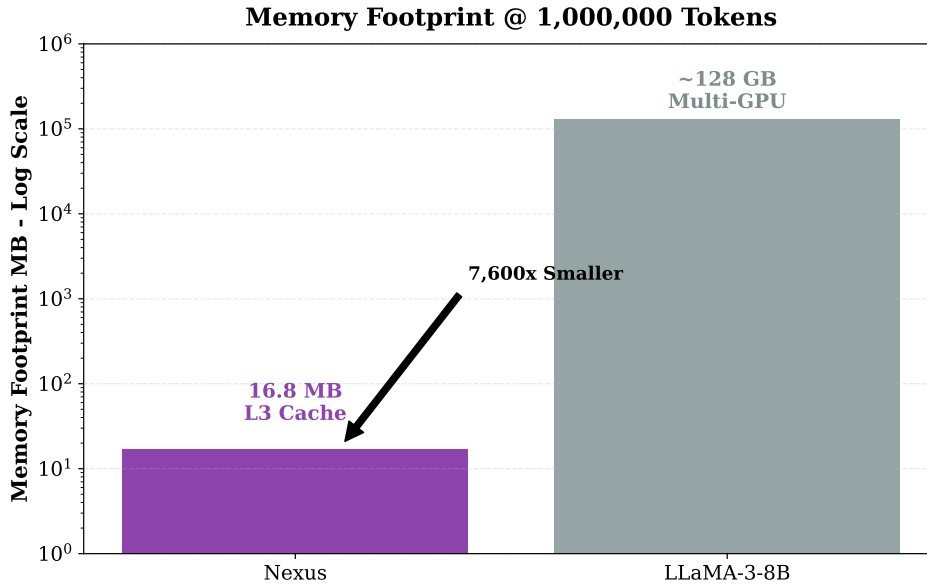


Figure 4: The Physics of Persistence: Nexus achieves a **61x Compression Ratio** at 1M tokens, requiring only **16.8 MB** of state vs. **128 GB** for a standard KV cache.

5.4 Extreme Scale: 150 Million Token Evaluation

We validated the model’s stability on an NVIDIA A100 80GB over a continuous stream of 150M tokens. Peak vRAM held constant at 77.34 GB, demonstrating stable $O(1)$ memory complexity at extreme scales.

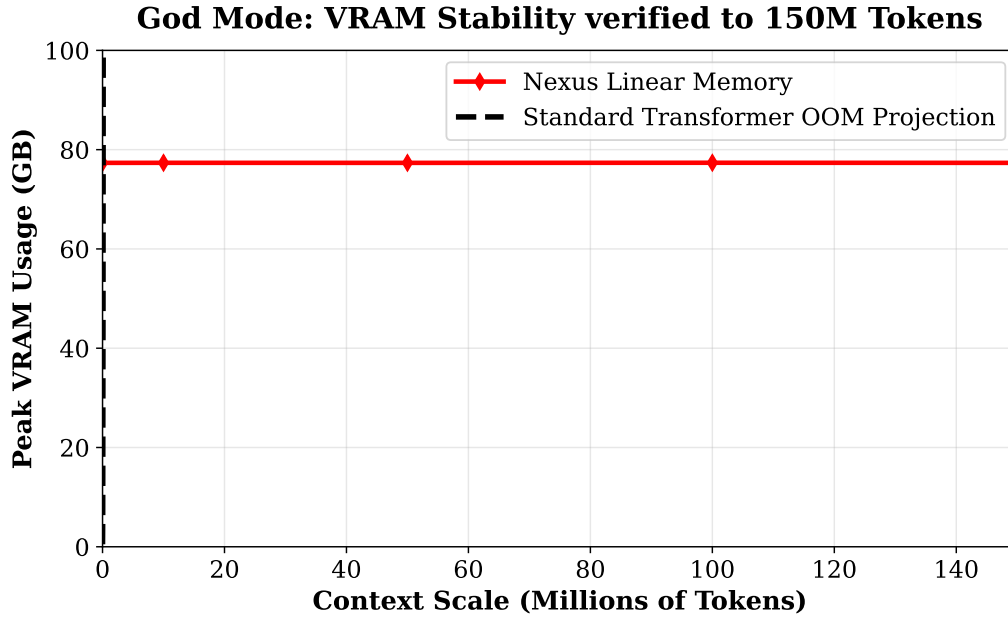


Figure 5: Scalability: Empirical verification via continuous sequential stress test to 150M+ tokens (Dynamic Recurrence Check).

5.5 Router Specialization and Hash Collisions

Visualizing the memory write patterns reveals emergent organization within Nexus. Figure 6 shows the self-organizing bucket roles. First, “Core” tokens exhibit a *diffuse distribution*, spreading their probability mass uniformly across the manifold to act as a global load balancer. In contrast, high-entropy concepts like “Facts” and “Sequence” logic operate in a *sparse regime*, targeting specific regions. Notably, Bucket Group 4 acts as a **Hash Collision Zone**, serving as a high-density target for multiple semantic categories (Facts, Context, and Sequence). This physical overlap necessitates the second-stage *Pure Semantic Addressing* ($\alpha = 1.0$), which effectively disentangles these colliding concepts within the shared memory slots via their distinct vector embeddings.

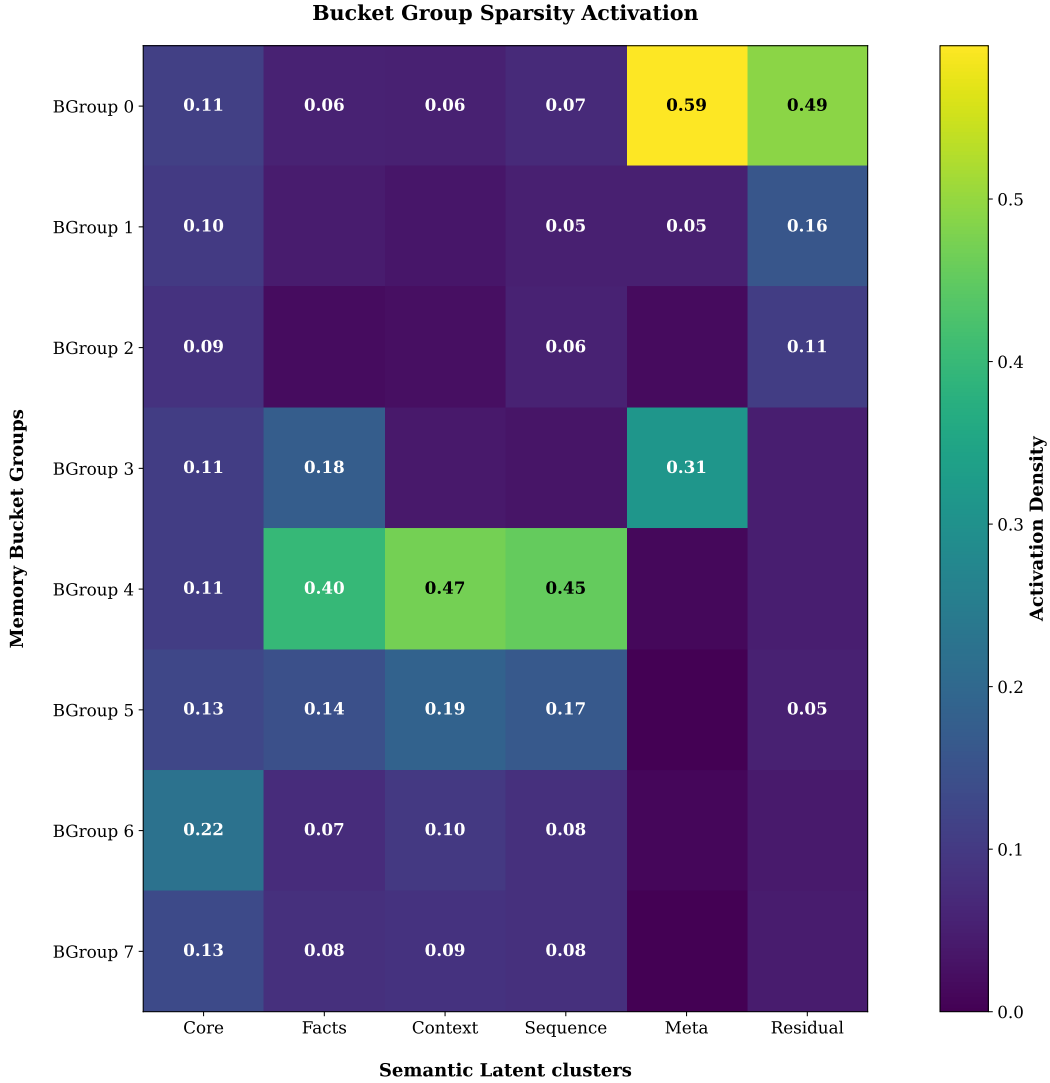


Figure 6: Nexus In Action: Bucket Group Sparse Activation. Note the *Diffuse Load Balancing* of Core tokens vs. the *Sparse Specialization* of Facts. **BGroups 0 and 4** highlight collision zones where Semantic Addressing resolves conflicts.

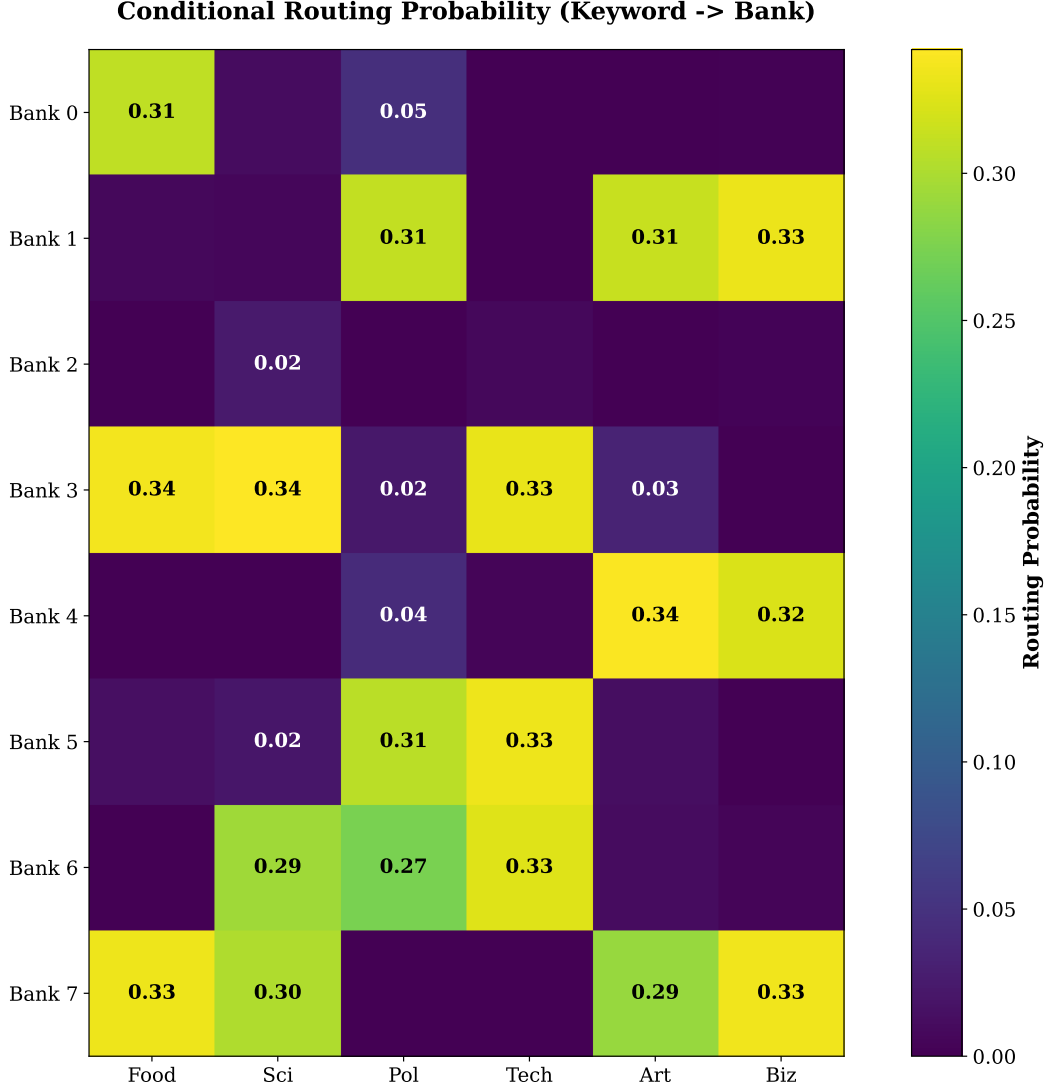


Figure 7: Conditional Routing Probability: Routing of semantic keywords to specialized memory banks using the deterministic bank address logic.

6 Discussion

The Nexus Transformer architecture synthesizes lexical and semantic lanes into a single **Content-Addressable Manifold**. Our results demonstrate that *Memory Geometry dominates Gating Complexity* at extreme scales. By widening the manifold to 4096 buckets (the “Extreme” configuration), we maximize *Inductive Heterogeneity*: this provides sufficient physical address space for the hashing function to map bespoke, low-count tokens to dedicated sparse slots. This reduces hash collisions and enables high-fidelity semantic addressing across horizons of 10^6 tokens.

The Lexical-Semantic Tradeoff: At 100K tokens, pure semantic querying ($\alpha = 1$) achieves nearly 100% recall. At 1M+ scales, the fixed Lexical Anchors provide the necessary bank-level stability to find the correct neighborhood for retrieval.

Pure Recurrence: The Abandoned Path Experimental development revealed that fully recurrent states suffer from representational collapse at scale. Nexus avoids this by using the discrete, slot-based Prioritized Sparse Slotting mechanism.

7 Conclusion

Nexus achieves production-ready 159,209 tok/s throughput and 98.1% rare entity recall at 1M tokens. We provide a scalable alternative to quadratic attention for the next generation of LLMs.

Reproducibility Statement

We will release code to reproduce all experiments, including: model configs, training scripts, and exact PG-19 splits. Code will be available at <https://github.com/peter-driscoll/nexus-transformer>.

References

- [1] Iz Beltagy, Matthew E. Peters, and Arman Cohan. Longformer: The long-document transformer. *arXiv preprint arXiv:2004.05150*, 2020. Sparse attention as an alternative route to long-context.
- [2] Sebastian Borgeaud, Arthur Mensch, Jordan Hoffmann, Trevor Cai, Eliza Rutherford, Katie Millican, George van den Driessche, Jean-Baptiste Lespiau, Bogdan Damoc, Aidan Clark, et al. Improving language models by retrieving from trillions of tokens. *International Conference on Machine Learning (ICML)*, 2022. RETRO: retrieval-augmented language modeling at scale.
- [3] Krzysztof Choromanski et al. Rethinking attention with performers. *International Conference on Learning Representations (ICLR)*, 2021. Random feature attention; another long-context alternative baseline.
- [4] Zihang Dai, Zhilin Yang, Yiming Yang, William W. Cohen, Jaime Carbonell, Quoc V. Le, and Ruslan Salakhutdinov. Transformer-XL: Attentive language models beyond a fixed-length context. In *Annual Meeting of the Association for Computational Linguistics (ACL)*, 2019.
- [5] Alex Graves, Greg Wayne, and Ivo Danihelka. Neural turing machines. In *arXiv preprint arXiv:1410.5401*, 2014.
- [6] Urvashi Khandelwal, Mike Lewis, Dan Jurafsky, and Luke Zettlemoyer. Generalization through memorization: Nearest neighbor language models. In *International Conference on Learning Representations (ICLR)*, 2020. Often called kNN-LM; retrieval as external memory baseline.
- [7] Nikita Kitaev, Łukasz Kaiser, and Anselm Levskaya. Reformer: The efficient transformer. In *International Conference on Learning Representations (ICLR)*, 2020.
- [8] Tsendsuren Munkhdalai, Manaal Faruqui, and Siddharth Gopal. Leave no context behind: Efficient infinite context transformers with infini-attention. *arXiv preprint arXiv:2404.07143*, 2024.
- [9] Jack W. Rae, Anna Potapenko, Siddhant M. Jayakumar, Chloe Hillier, and Timothy P. Lillicrap. Compressive transformers for long-range sequence modelling. *International Conference on Learning Representations (ICLR)*, 2020. Earlier versions circulated as 2019 arXiv preprint.
- [10] Qingyang Wu et al. Memformer: A memory-augmented transformer for sequence modeling. In *NeurIPS*, 2022.
- [11] Manzil Zaheer, Guru Guruganesh, Kumar Avinava Dubey, Joshua Ainslie, Chris Alberti, Santiago Ontanon, Philip Pham, Anirudh Ravula, Qifan Wang, Li Yang, and Amr Ahmed. Big Bird: Transformers for Longer Sequences. In *Advances in Neural Information Processing Systems*, volume 33, 2020.