

# RAM-NET: EXPRESSIVE LINEAR ATTENTION WITH SELECTIVELY ADDRESSABLE MEMORY

**Kaicheng Xiao, Haotian Li, Liran Dong, Guoliang Xing**

The Chinese University of Hong Kong

{xk023, lh023, dl123, glxing}@ie.cuhk.edu.hk

## ABSTRACT

While linear attention architectures offer efficient inference, compressing unbounded history into a fixed-size memory inherently limits expressivity and causes information loss. To address this limitation, we introduce Random Access Memory Network (RAM-Net), a novel architecture designed to bridge the gap between the representational capacity of full attention and the memory efficiency of linear models. The core of RAM-Net maps inputs to high-dimensional sparse vectors serving as explicit addresses, allowing the model to selectively access a massive memory state. This design enables exponential state size scaling without additional parameters, which significantly mitigates signal interference and enhances retrieval fidelity. Moreover, the inherent sparsity ensures exceptional computational efficiency, as state updates are confined to minimal entries. Extensive experiments demonstrate that RAM-Net consistently surpasses state-of-the-art baselines in fine-grained long-range retrieval tasks and achieves competitive performance in standard language modeling and zero-shot commonsense reasoning benchmarks, validating its superior capability to capture complex dependencies with significantly reduced computational overhead.

## 1 INTRODUCTION

Transformer has emerged as the dominant architecture for sequence modeling due to their effectiveness in capturing long-range dependencies (Vaswani et al., 2017). However, despite their widespread success, they fundamentally incur a linear memory growth during inference. The standard full attention mechanism necessitates caching the complete history of key-value pairs to ensure retrieval precision, creating a substantial memory bottleneck for long sequences (Dao et al., 2022). Moreover, standard transformers lack an explicit mechanism to selectively discard obsolete information (Rae et al., 2019). Instead, they simply concatenate new tokens to the existing historical context, treating every input as a permanent entry in the memory. Consequently, this passive accumulation of data leads to significant computational inefficiencies and attention dilution when processing extensive contexts.

To mitigate the quadratic computational cost and linear memory growth of full attention, substantial research has been dedicated to linearizing the attention mechanism. Existing approaches range from kernel-based approximations (Katharopoulos et al., 2020; Choromanski et al., 2020; Wang et al., 2020) to recurrent architectures and state space models (Gu et al., 2021; Poli et al., 2023; Gu & Dao, 2024; Peng et al., 2023; Ma et al., 2022; Qin et al., 2024). Despite their diversity, these methods fundamentally rely on compressing unbounded historical context into a fixed-size recurrent state. While this formulation successfully achieves constant memory and linear time complexity during inference, it introduces a significant expressivity bottleneck. Specifically, compressing the entire history into a single fixed-size state forces distinct contextual features to compete for limited storage (Sun et al., 2024). Unlike full attention, which retains direct access to specific past key-value pairs, linear architectures rely on this superposed state for retrieval. Consequently, this compact representation often compromises the ability to model fine-grained long-range interactions, particularly in tasks where high-resolution recall is essential (Arora et al., 2024b).

To bridge the gap between representational capacity and memory efficiency, we introduce RAM-Net. This architecture is designed to reconcile the high-fidelity retrieval capabilities of full attention

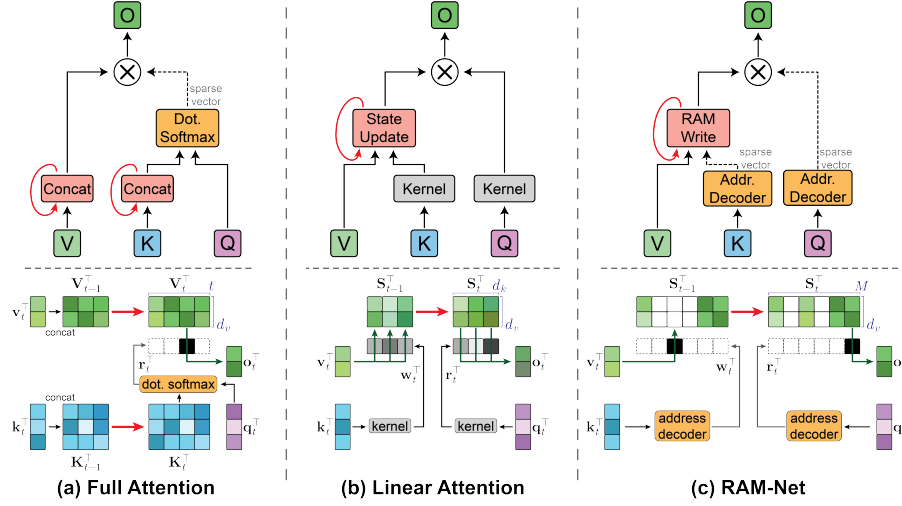


Figure 1: Comparison of memory mechanisms. (a) Full Attention: Retains the entire history for retrieval, resulting in linear memory growth. (b) Linear Attention: Compresses history into a fixed-size state. The reliance on kernel-based similarity leads to limited capacity and inevitable interference. (c) RAM-Net: Decouples memory capacity from feature dimension via an Address Decoder, which maps dense vectors  $\mathbf{k}_t$  and  $\mathbf{v}_t$  into high-dimensional sparse addresses  $\mathbf{w}_t$  and  $\mathbf{r}_t$ . This enables massive state capacity and high-fidelity retrieval with constant memory state size.

with the constant memory overhead of linear models. Departing from the traditional paradigm of compressing history into a fixed-size latent bottleneck, the core of RAM-Net utilizes a differentiable address decoder. This mechanism projects dense low-dimensional input vectors into high-dimensional sparse vectors that serve as explicit addresses, enabling the model to selectively access a massive global memory state. This design allows for exponential state size scaling without any increase in learnable parameters. By distributing historical context across this expanded memory space, RAM-Net effectively isolates diverse semantic features, significantly mitigating the signal interference and information loss inevitably caused by compact state compression. Moreover, despite the massive memory capacity, RAM-Net maintains exceptional computational efficiency. The inherent sparsity ensures that state updates and retrieval operations are strictly confined to a minimal set of active entries. This mechanism mimics the precision of Random Access Memory, enabling the model to capture fine-grained, long-range dependencies with high precision while keeping computational costs low.

## 2 BACKGROUND

To provide a unified view of sequence modeling, we revisit the attention mechanism from the perspective of state maintenance. At each step  $t$ , the model maintains a state  $\mathbf{S}_t$  by updating and retrieving historical information from it.

Given an input row vector  $\mathbf{x}_t \in \mathbb{R}^{1 \times d_m}$ , we first obtain the query  $\mathbf{q}_t, \mathbf{k}_t \in \mathbb{R}^{1 \times d_k}$  and value  $\mathbf{v}_t \in \mathbb{R}^{1 \times d_v}$  via learnable projections  $\mathbf{W}_q, \mathbf{W}_k \in \mathbb{R}^{d_m \times d_k}$  and  $\mathbf{W}_v \in \mathbb{R}^{d_m \times d_v}$ . In this framework, we interpret the attention process as two distinct operations: the key and value serve to update the memory state via a *write* operator  $\mathcal{W}$ , while the query performs information retrieval via a *read* operator  $\mathcal{R}$ :

$$\mathbf{S}_t = \mathcal{W}(\mathbf{S}_{t-1}, \mathbf{k}_t, \mathbf{v}_t), \quad \mathbf{o}_t = \mathcal{R}(\mathbf{S}_t, \mathbf{q}_t). \quad (1)$$

The output is subsequently projected by  $\mathbf{W}_o$ . Different attention architectures fundamentally differ in how they structure the memory state  $\mathbf{S}_t$  and define these specific update ( $\mathcal{W}$ ) and retrieval ( $\mathcal{R}$ ) mechanisms (Arora et al., 2024b).

## 2.1 FULL ATTENTION

Full Attention (Vaswani et al., 2017) explicitly maintains the entire history as its memory state. As illustrated in Fig. 1 (a), the state at step  $t$  is defined as a tuple  $\mathbf{S}_t = (\mathbf{K}_t, \mathbf{V}_t)$ , consisting of the accumulated key matrix  $\mathbf{K}_t \in \mathbb{R}^{t \times d_k}$  and value matrix  $\mathbf{V}_t \in \mathbb{R}^{t \times d_v}$ .

Specifically, the update operator  $\mathcal{W}$  appends the current inputs to the history cache, while the read operator  $\mathcal{R}$  computes Softmax-normalized scores  $\mathbf{r}_t \in \mathbb{R}^{1 \times t}$  to retrieve context via the query-key matching:

$$\begin{aligned}\mathbf{K}_t &= \text{Concat}(\mathbf{K}_{t-1}, \mathbf{k}_t), \\ \mathbf{V}_t &= \text{Concat}(\mathbf{V}_{t-1}, \mathbf{v}_t), \\ \mathbf{r}_t &= \text{Softmax}(\mathbf{q}_t \mathbf{K}_t^\top), \\ \mathbf{o}_t &= \mathbf{r}_t \mathbf{V}_t.\end{aligned}\tag{2}$$

However, since the size of state  $\mathbf{S}_t$  grows linearly with the sequence length  $t$ , this approach incurs significant memory and computational costs during inference.

## 2.2 LINEAR ATTENTION

Linear attention restricts the state to a fixed-size  $\mathbf{S}_t \in \mathbb{R}^{d_k \times d_v}$ , as shown in Fig. 1 (b). To compress unbounded history into these finite state, linear attention introduces a kernel function  $\phi(\cdot)$ . It uses the mapped dense vectors  $\mathbf{w}_t = \phi(\mathbf{k}_t) \in \mathbb{R}^{1 \times d_k}$  and  $\mathbf{r}_t = \phi(\mathbf{q}_t) \in \mathbb{R}^{1 \times d_k}$  to explicitly control the writing and reading of the state. The state update  $\mathcal{W}$  becomes a recurrent accumulation, and retrieval  $\mathcal{R}$  becomes a linear projection:

$$\mathbf{S}_t = g(\mathbf{x}_t) \mathbf{S}_{t-1} + \mathbf{w}_t^\top \mathbf{v}_t, \quad \mathbf{o}_t \propto \mathbf{r}_t \mathbf{S}_t.\tag{3}$$

Here,  $g(\mathbf{x}_t)$  acts as a decay gate to alleviate capacity exhaustion via selective forgetting, a mechanism widely used in modern variants.

However, the memory capacity is tightly coupled with the feature dimension  $d_k$ , which defines the number of memory rows. Scaling up capacity to reduce information overlap necessitates increasing  $d_k$ , which causes the memory and computational cost of the state update to grow synchronously. This coupling makes it inefficient to expand the state for complex contexts, leading to the superposition of memories (Sun et al., 2024) and inaccurate retrieval.

## 3 RAM-NET

To scale the memory capacity without incurring additional parameter or computational overheads, we introduce **RAM-Net**, as illustrated in Fig. 1 (c). The core design is an efficient mapping mechanism called *address decoding* that transforms the dense, low-dimensional vectors  $\mathbf{k}_t, \mathbf{q}_t \in \mathbb{R}^{1 \times d_k}$  into high-dimensional sparse write and read vectors  $\mathbf{w}_t, \mathbf{r}_t \in \mathbb{R}^{1 \times M}$  (where  $M \gg d_k$ ). Specifically, we enforce Top- $K$  sparsity on  $\mathbf{w}_t$  and  $\mathbf{r}_t$ , resulting in  $K$ -hot vectors that act as explicit addresses. Correspondingly, we define the memory state  $\mathbf{S}_t \in \mathbb{R}^{M \times d_v}$  to comprise  $M$  discrete memory slots. In this framework, the non-zero indices of  $\mathbf{w}_t$  and  $\mathbf{r}_t$  act as pointers, activating only  $K$  specific slots for each read or write operation—a sparse access pattern analogous to Random Access Memory (RAM).

This architecture yields two critical advantages. First, decoupling the memory capacity  $M$  from the feature dimension  $d_k$  enables massive state size scaling without inflating model parameters. The resulting large capacity of  $\mathbf{S}_t$  significantly reduces signal interference, ensuring high-fidelity retrieval. Second, since only  $K$  of the  $M$  memory slots are active at each write or read operation, the model achieves a highly efficient computational cost of  $\mathcal{O}(K \cdot d_v)$ .

### 3.1 ADDRESS DECODER

To derive high-dimensional sparse addresses, we propose an efficient address decoding method  $\mathcal{A}_{K,U}(\cdot)$ , which is composed of a *Product Softmax* expansion  $\rho_U(\cdot)$  that expands the input into a concentrated high-dimensional distribution, and a Top- $K$  truncation  $\mathcal{T}_K(\cdot)$  that enforces sparsity by discarding minor values.

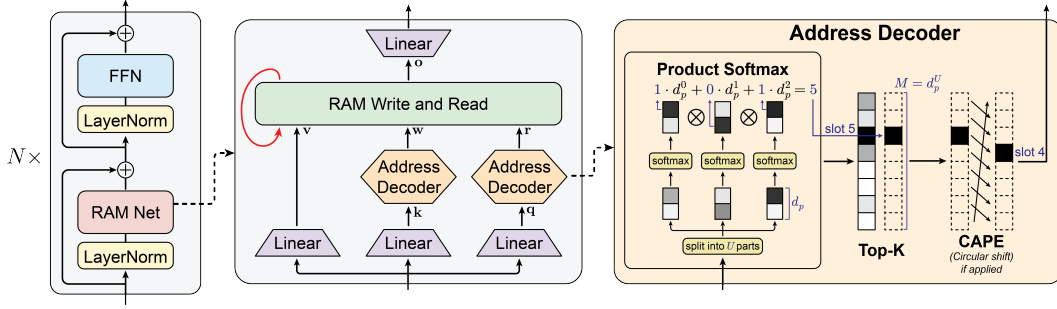


Figure 2: Overview of the RAM-Net architecture. The Address Decoder transforms  $\mathbf{k}_t$  and  $\mathbf{q}_t$  vectors into high-dimensional sparse addresses via Product Softmax, Top- $K$  truncation, and Cyclic Address Positional Embedding (CAPE). For visual clarity, we illustrate a simplified configuration with  $U = 3$  partitions and sub-dimension  $d_p = 2$ . This results in a total memory capacity of  $M = 8$  slots with selection sparsity  $K = 1$ .

**$U$ -order product softmax  $\rho_U(\cdot)$ .** As shown in Fig. 2, we partition the vector  $\mathbf{k}_t$  (and similarly  $\mathbf{q}_t$ ) into  $U$  independent sub-vectors  $\mathbf{k}_t = [\mathbf{k}_t^{(1)}, \dots, \mathbf{k}_t^{(U)}]$ , where each component has a dimension of  $d_p = d_k/U$ . We then synthesize the full address vector over  $M = (d_p)^U$  slots via the Kronecker product ( $\otimes$ ) with temperature scaling parameter  $\tau$ , defined as the  $U$ -order product softmax (where  $U$  denotes the construction order)::

$$\rho_U(\mathbf{k}_t) = \bigotimes_{u=1}^U \text{Softmax} \left( \frac{\mathbf{k}_t^{(u)}}{\tau} \right). \quad (4)$$

Notably, this enables scaling the memory capacity of state  $\mathbf{S}_t$  to the high-dimensional space  $M$  without introducing any additional model parameters. Moreover, compared to directly applying Softmax over the massive dimension  $M$ , this product formulation ensures more effective gradient propagation and optimization stability, leading to superior performance (see detailed analysis in Sec. 5.3).

**Top- $K$  truncation  $\mathcal{T}_K(\cdot)$ .** Subsequently, the truncation operator  $\mathcal{T}_K$  enforces explicit sparsity by retaining only the  $K$  largest values and zeroing out the rest. The complete formulation of the address decoding is given by:

$$\mathcal{A}_{K,U}(\mathbf{k}_t) = \mathcal{T}_K(\rho_U(\mathbf{k}_t)). \quad (5)$$

This sparse addressing significantly reduces the computational complexity of subsequent state updates and retrieval operations, as only  $K$  active slots require processing per step. At the same time, top- $K$  can be combined with the structure of product softmax to achieve fast address decoding (see detailed analysis in Sec. 4).

### 3.2 CYCLIC ADDRESS POSITIONAL EMBEDDING

Since the vectors  $\mathbf{q}_t$  and  $\mathbf{k}_t$  are derived solely from the input vector, the address decoding inherently lacks sequence order information. To bridge this gap, we draw inspiration from Rotary Positional Embeddings (RoPE) (Su et al., 2024) and propose *Cyclic Address Positional Embedding (CAPE)*, which incorporates relative positional information via a cyclic shift operator  $\mathcal{P}_t(\cdot)$ .

Specifically,  $\mathcal{P}_t$  performs a time-dependent circular shift on the address vector. For any vector  $\mathbf{a} \in \mathbb{R}^M$ , the operator acts as a cyclic permutation defined by  $[\mathcal{P}_t(\mathbf{a})]_i = \mathbf{a}_{(i+t) \bmod M}$ . By applying this transformation, the interaction between a write operation at step  $t$  and a read operation at step  $t'$  becomes dependent on the relative distance  $t - t'$ . This effectively converts absolute addressing into relative addressing, allowing the model to capture temporal structures analogous to a temporal convolution (especially when  $\mathbf{k}_t$  and  $\mathbf{q}_t$  are constant). In summary, the final write and read vectors  $\mathbf{w}_t$  and  $\mathbf{r}_t$  are formally defined as:

$$\begin{aligned} \mathbf{w}_t &= \mathcal{P}_t(\mathcal{A}_{K,U}(\mathbf{k}_t)), \\ \mathbf{r}_t &= \mathcal{P}_t(\mathcal{A}_{K,U}(\mathbf{q}_t)). \end{aligned} \quad (6)$$

### 3.3 MEMORY WRITE AND READ

Standard linear gating restricts memory updates to a rigid Exponential Moving Average (EMA), where the retention weight is determined by the complement of the write weight  $1 - \mathbf{w}_t^\top$ . This forces the model to forget historical information in exact proportion to the new input intensity, thereby preventing independent control over information preservation.

To address this, we introduce Power Decay Moving Average (PDMA), which generalizes the aggregation mechanism by decoupling the forgetting rate from the writing intensity. The update rules are defined as:

$$\begin{aligned} \mathbf{S}_t &= \text{Diag}(1 - \mathbf{w}_t)^\gamma \cdot \mathbf{S}_{t-1} + \mathbf{w}_t^\top \mathbf{v}_t, & \mathbf{S}_0 &= \mathbf{0}, \\ \mathbf{z}_t &= \text{Diag}(1 - \mathbf{w}_t)^\gamma \cdot \mathbf{z}_{t-1} + \mathbf{w}_t^\top, & \mathbf{z}_0 &= \mathbf{1}^\top / M, \\ \mathbf{o}_t &= \mathbf{r}_t (\text{Diag}(\mathbf{z}_t + \epsilon)^{-1} \cdot \mathbf{S}_t), \end{aligned} \quad (7)$$

where  $\mathbf{1} \in \mathbb{R}^{1 \times M}$  is a row vector of ones. Here, the hyperparameter  $\gamma \geq 0$  controls the non-linear decay term  $(1 - \mathbf{w}_t^\top)^\gamma$ , while  $\mathbf{z}_t$  tracks the accumulated weight mass for dynamic normalization. This formulation creates a continuous spectrum of memory behaviors: as  $\gamma \rightarrow 0$ , the decay term approaches unity, transforming the update into a Weighted Cumulative Mean that preserves long-term dependencies without attenuation. Conversely, setting  $\gamma > 1$  accelerates forgetting beyond the standard linear decay ( $\gamma = 1$ ), enabling the model to aggressively suppress historical noise in favor of the most recent context.

## 4 IMPLEMENTATION

**Architecture configuration.** We employ a hybrid strategy that alternates between relative and absolute addressing modes. Shallow layers primarily utilize positional embedding to capture local temporal dynamics. In contrast, deeper layers primarily employ absolute addressing by omitting the cyclic shift operator  $\mathcal{P}_t$  in equ. 6.

**Training convergence acceleration.** We introduce a dynamic scalar re-parameterization. For the projection matrices  $\mathbf{W}_q$  and  $\mathbf{W}_k$ , we decompose each weight  $w_{ij}$  into  $e^\alpha \cdot w'_{ij}$ , where  $\alpha$  is a learnable per-head parameter. This learnable scalar is equivalent to a dynamic temperature of softmax, ensuring effective gradient flow and facilitating faster convergence.

**Gradient stability.** For the non-linear decay term  $(1 - \mathbf{w}_t^\top)^\gamma$  in PDMA, we address the critical gradient issues near the boundary  $\mathbf{w}_t \rightarrow \mathbf{1}$ , where the derivative becomes infinity when  $\gamma < 1$ , and vanishes to zero when  $\gamma \geq 1$ . To resolve this, we employ a proxy gradient technique during backpropagation, substituting the original derivative with that of a smoothed function  $(\epsilon + (1 - \epsilon)(1 - \mathbf{w}_t^\top))^\gamma$ . This separation of computation paths ensures robust gradient flow without compromising the exactness of the forward calculation.

**Computational efficiency.** During address decoding, directly applying Top- $K$  selection on an address vector of size  $M = d_p^U$  entails prohibitive computational costs. To circumvent this, we avoid explicit vector construction by leveraging the combinatorial structure of the addresses. Specifically, we employ a log-domain beam search strategy that iteratively merges and sorted candidates. This approach ensures numerical stability and significantly reduces retrieval complexity from  $\mathcal{O}(d_p^U)$  to  $\mathcal{O}(U \cdot K^2 + U \cdot d_p \log d_p)$ .

To optimize execution, we develop specialized kernels for both training and inference. During parallel training, we implement a segment-based approach where sequences are sorted to aggregate sparse read/write operations into contiguous time-slots, enabling efficient batched computation. For autoregressive inference, we leverage sparsity to directly access memory slots via decoded indices, retrieving values without full-state iteration.

## 5 EXPERIMENTS

### 5.1 SYNTHETIC BENCHMARKS

We evaluate the retrieval capability of our model using the multi-query associative recall (MQAR) task from the Zoology framework (Arora et al., 2024a). This task requires the model to memorize

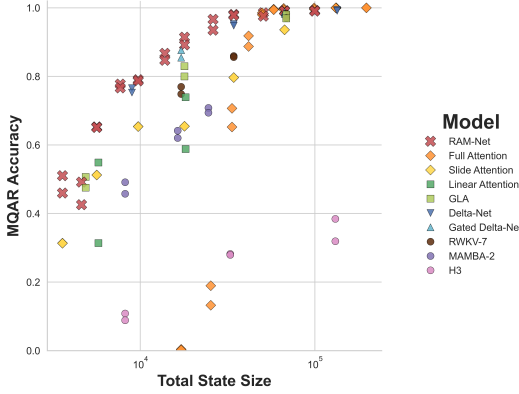


Figure 3: MQAR accuracy vs. total state size.

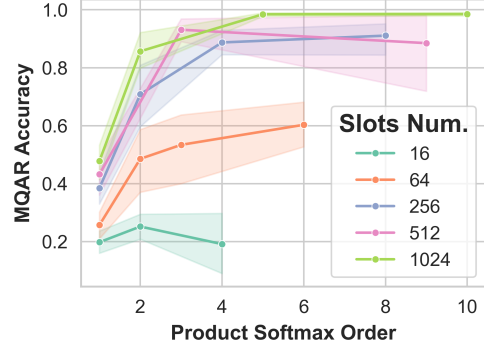


Figure 4: Ablation study of product softmax order  $U$ .

a sequence of key-value pairs embedded within a long context and subsequently retrieve the correct value associated with a specific query key. Our main baselines include full attention (Vaswani et al., 2017), linear attention (Katharopoulos et al., 2020), sliding window attention (Beltagy et al., 2020), GLA (Yang et al., 2024), DeltaNet (Schlag et al., 2021a), Gated DeltaNet (Yang et al., 2024), RWKV-7 (Peng et al., 2025), Mamba2 (Dao & Gu, 2024), and H3 (Fu et al., 2022). We test these models across various configurations to cover a spectrum of memory capacities. All reported results are obtained from our own implementations.

We present the comparative results in Fig. 3, which plots the MQAR accuracy against the state size, representing the total memory overhead required during inference. Following the default Zoology configuration, we report the uniform average accuracy across seven settings denoted as  $(N_{\text{pairs}}, L_{\text{seq}})$ , representing the number of pairs and sequence length: (4, 64), (8, 64), (16, 64), (32, 128), (64, 256), (128, 512), and (256, 1024). As illustrated, RAM-Net consistently achieves superior retrieval accuracy across various state sizes. This effectively validates that our memory mechanism demonstrates superior information storage efficiency compared to other architectures.

## 5.2 LANGUAGE MODELING

**Baselines and Experimental Setup.** We compare our method against baselines including Transformer++ (Touvron et al., 2023), GLA (Yang et al., 2024), HGRN2 (Qin et al., 2024), Gated DeltaNet (Yang et al., 2024), and Mamba2 (Dao & Gu, 2024). To ensure a fair comparison, all models are pre-trained under a unified configuration using the Flame framework (Zhang & Yang, 2025), utilizing implementations from Flash Linear Attention (Yang & Zhang, 2024). We fix the hidden size (model width) to 1024, and scale to 340M. Training is conducted on the FineWeb-Edu dataset (Lozhkov et al., 2024) with a budget of 10B tokens. We employ the Llama 2 tokenizer (32k vocabulary) and a context window of 4,096 tokens. We optimize the model using AdamW with a peak learning rate of  $1.0 \times 10^{-3}$ , a weight decay of 0.1, and gradient clipping of 1.0. The learning rate follows a cosine decay schedule with a 500M token warm-up. The global batch size is set to 0.5M tokens, distributed across 8 NVIDIA H200 GPUs.

**Evaluation.** We evaluate performance using the Language Model Evaluation Harness (Gao et al., 2024) on standard benchmarks: WikiText-103 (Merity et al., 2016), MMLU (Hendrycks et al., 2021), ARC-Challenge/Easy (Clark et al., 2018), OpenbookQA (Mihaylov et al., 2018), SciQ (Welbl et al., 2017), COPA (Roemmele et al., 2011), PIQA (Bisk et al., 2020), HellaSwag (Zellers et al., 2019), and WinoGrande (Sakaguchi et al., 2019). The comparative results are summarized in Table 1. RAM-Net achieves comparable performance with other SOTA methods.

Table 1: Zero-shot performance comparison on language modeling and common-sense reasoning.

Model	Wiki. ppl ↓	MMLU acc ↑	ARC-e acc ↑	ARC-c acc ↑	OBQA acc ↑	SciQ acc ↑	COPA acc ↑	PIQA acc ↑	Hella. acc ↑	Wino. acc ↑	Params. <sup>2</sup> size ↓	Active State <sup>3</sup> per token ↓
Transformer++	35.96	23.0	55.7	24.0	<b>34.2</b>	82.9	67.0	65.0	32.5	50.6	341.1M	50.3M
GLA	34.37	23.0	56.1	21.3	31.2	79.1	66.0	64.1	31.5	50.8	341.7M	<u>3.1M</u>
HGRN2	30.58	23.2	<b>59.1</b>	<b>24.6</b>	33.0	80.7	<b>73.0</b>	<b>67.2</b>	32.9	51.1	<u>341.1M</u>	<u>3.1M</u>
Gated DeltaNet <sup>1</sup>	<b>28.79</b>	23.0	<u>59.0</u>	23.5	32.2	<b>83.1</b>	66.0	66.7	<b>33.2</b>	<u>51.7</u>	347.0M	8.5M
Mamba2 <sup>1</sup>	<u>29.96</u>	<b>24.1</b>	56.8	24.4	32.2	81.4	<u>71.0</u>	<u>66.8</u>	<u>33.1</u>	<b>51.9</b>	349.6M	12.9M
RAM-Net (Top-8)	32.33	<u>23.3</u>	57.1	<b>24.6</b>	<b>34.2</b>	79.4	68.0	65.3	31.7	50.9	<b>340.8M</b>	<b>0.4M</b>

<sup>1</sup> Include an additional cross-token short convolution layer with kernel size  $k = 4$ ; other models use token-wise projections only.

<sup>2</sup> The number of trainable parameters, while the word embedding parameters are excluded to ensure consistency with tied word embeddings.

<sup>3</sup> Number of activated states update/read per token, directly reflecting compute and memory-bandwidth demand (Transformer++ is calculated at sequence length  $L = 1024$ ).

### 5.3 ABLATION STUDY

We conduct an ablation study to validate the design choices of the Address Decoder. Specifically, we decouple the impact of product softmax order  $U$  (the Product Softmax order of  $\mathbf{k}$  and  $\mathbf{q}$ ) from the total memory capacity  $M$ , analyzing how each factor contributes to the model’s retrieval capability.

**Product Softmax Order  $U$ .** First, we analyze the impact of the Product Softmax order  $U$ . We treat the setting  $U = 1$  as a baseline, where the network applies softmax directly across the entire massive dimension of memory slots without decomposition. In contrast, configurations with  $U > 1$  employ our proposed product softmax formulation to derive addresses from smaller sub-vectors.

As illustrated in Fig. 4, we observe a significant improvement in MQAR accuracy as  $U$  increases. We attribute this performance gain to the superior gradient dynamics facilitated by the product softmax structure. For the non-decomposed baseline ( $U = 1$ ), the combination of global Softmax and Top- $K$  truncation results in extreme gradient sparsity on  $\mathbf{r}$  and  $\mathbf{w}$ , where at most  $K/M = K/(d_p)^U$  elements receive non-zero gradients. This sparsity induces high variance in gradient estimation, leading to optimization instability. In contrast, the Kronecker product structure acts as an efficient gradient distributor. During backpropagation, it distributes effective gradients from  $\mathbf{r}$  and  $\mathbf{w}$  to all  $U$  sub-vectors, ensuring that each sub-vector receives gradient feedback on at least  $1/d_p$  of its elements. This mechanism effectively reduces gradient variance, ensuring stable convergence even with high-dimensional memory states.

**Memory Capacity  $M$ .** Next, we examine the scalability of memory capacity  $M$ . The results demonstrate a clear positive correlation between the number of memory slots and retrieval performance. As shown in Fig. 4, increasing  $M$  consistently leads to higher MQAR accuracy. A larger memory capacity enables the model to map distinct semantic features to non-overlapping slots more effectively. This reduction in superposition preserves high-fidelity information and ensures higher precision in long-range retrieval.

## 6 RELATED WORK

**Linear Sequence Modeling** To mitigate the quadratic complexity of standard Transformers, linear sequence models compress the unbounded historical context into a fixed-size recurrent state (Katharopoulos et al., 2020; Sun et al., 2023). Early approaches, such as Linear Transformers (Katharopoulos et al., 2020) and RetNet (Sun et al., 2023), typically employ kernel-based feature maps or time-invariant decay mechanisms to aggregate context. While efficient, these methods rely on static dynamics that fail to adaptively filter information based on input content.

To enhance expressivity, recent architectures have introduced data-dependent dynamics. Models like Mamba2 (Dao & Gu, 2024), RWKV7 (Peng et al., 2023), and HGRN2 (Qin et al., 2024) utilize input-driven gating mechanisms to actively modulate information retention. Furthermore, approaches such as DeltaNet (Schlag et al., 2021b), Gated DeltaNet (Yang et al., 2024), and Kimi Linear (Team et al., 2025) adopt the Delta Rule, formulating the state update as a dynamic retrieve-and-rewrite process. This mechanism can be generalized as  $\mathbf{S}_t = \mathcal{W}(\mathbf{S}_{t-1}, \mathbf{k}_t, \beta(\mathbf{v}_t - \mathcal{R}(\mathbf{S}_{t-1}, \mathbf{k}_t)))$ , where  $\mathcal{R}$  represents the retrieval of existing knowledge from the state  $\mathbf{S}_{t-1}$  using the key  $\mathbf{k}_t$ , and the update operator  $\mathcal{W}$  incorporates the error-driven correction weighted by  $\beta$ . This formulation

interprets the recurrent update as an online optimization step, allowing for more precise memory editing (Schlag et al., 2021b; Sun et al., 2024).

Despite these advancements in update rules, purely linear models share a fundamental bottleneck: the fixed state capacity. Compressing a long sequence into a compact hidden state inevitably results in memory collisions and information loss, limiting performance on fine-grained retrieval tasks over long horizons. In contrast, RAM-Net projects historical context into an expanded high-dimensional memory space. By leveraging sparse read-write mechanisms, our approach isolates distinct semantic features into separate slots, thereby minimizing memory superposition and ensuring high-fidelity retrieval.

**Vector Quantization and Explicit Memory** To bypass the limitations of compressed recurrent states, several architectures employ explicit memory slots as storage units. Approaches such as Transformer-VQ (Lingle, 2023) and PQCache (Zhang et al., 2025) utilize Vector Quantization (VQ) to map inputs to discrete memory entries. While efficient, these methods rely on codebook-based addressing involving non-differentiable selection operators (e.g.,  $\text{argmax}$ ), which necessitates approximation techniques like the Straight-Through Estimator (STE) (Bengio et al., 2013) for training.

In contrast to discrete quantization, Neural Turing Machines (NTM) (Graves et al., 2014) employ differentiable memory mechanisms. However, they rely on simulating the movement of a read-write head to access memory, resulting in an indirect slot selection mechanism. Conversely, RAM-Net distinguishes itself by implementing a fully differentiable address decoder that projects inputs into sparse high-dimensional address vectors, enabling direct slot selection. This design avoids the non-differentiable operations present in VQ methods while eliminating the sequential head-shifting logic of NTMs.

## 7 DISCUSSION

The explicit addressing design of RAM-Net provides inherent properties beneficial for both model interpretability and system efficiency. By explicitly tracking memory updates and retrievals via vectors  $\mathbf{w}_t$  and  $\mathbf{r}_t$ , the model offers granular insights into token interactions and the functional roles of individual memory slots. Moreover, the sparse access pattern is particularly advantageous for hardware-constrained environments. It allows for hierarchical memory management, where the full state resides in cheaper storage (e.g., CPU RAM or SSD), and only frequently accessed hot slots are dynamically cached in VRAM via policies like Least Recently Used (LRU), significantly reducing the GPU memory overhead.

## 8 CONCLUSION

In this work, we propose RAM-Net to reconcile the trade-off between the high fidelity of full attention and the memory efficiency of linear models. RAM-Net introduces a differentiable address decoding mechanism that maps inputs to explicit high-dimensional memory slots. This paradigm allows state capacity to scale independently of model parameters. Empirical validation across retrieval and language modeling tasks demonstrates that RAM-Net offers a robust solution for processing complex sequences while maintaining the efficiency benefits of sparse computation.

## REFERENCES

- Simran Arora, Sabri Eyuboglu, Aman Timalsina, Isys Johnson, Michael Poli, James Zou, Atri Rudra, and Christopher Ré. Zoology: Measuring and improving recall in efficient language models. In *Proceedings of 12th International Conference on Learning Representations (ICLR)*. ICLR, 2024a.
- Simran Arora, Sabri Eyuboglu, Michael Zhang, Aman Timalsina, Silas Alberti, Dylan Zinsley, James Zou, Atri Rudra, and Christopher Ré. Simple linear attention language models balance the recall-throughput tradeoff. *arXiv preprint arXiv:2402.18668*, 2024b.
- Iz Beltagy, Matthew E Peters, and Arman Cohan. Longformer: The long-document transformer. *arXiv preprint arXiv:2004.05150*, 2020.



- Yoshua Bengio, Nicholas Léonard, and Aaron Courville. Estimating or propagating gradients through stochastic neurons for conditional computation. *arXiv preprint arXiv:1308.3432*, 2013.
- Yonatan Bisk, Rowan Zellers, Jianfeng Gao, Yejin Choi, et al. Piqa: Reasoning about physical commonsense in natural language. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pp. 7432–7439, 2020.
- Krzysztof Choromanski, Valerii Likhoshesterov, David Dohan, Xingyou Song, Andreea Gane, Tamas Sarlos, Peter Hawkins, Jared Davis, Afroz Mohiuddin, Lukasz Kaiser, et al. Rethinking attention with performers. *arXiv preprint arXiv:2009.14794*, 2020.
- Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge, 2018. URL <https://arxiv.org/abs/1803.05457>.
- Tri Dao and Albert Gu. Transformers are ssms: Generalized models and efficient algorithms through structured state space duality. *arXiv preprint arXiv:2405.21060*, 2024.
- Tri Dao, Dan Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. Flashattention: Fast and memory-efficient exact attention with io-awareness. *Advances in neural information processing systems*, 35:16344–16359, 2022.
- Daniel Y Fu, Tri Dao, Khaled K Saab, Armin W Thomas, Atri Rudra, and Christopher Ré. Hungry hungry hippos: Towards language modeling with state space models. *arXiv preprint arXiv:2212.14052*, 2022.
- Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Alain Le Noac’h, Haonan Li, Kyle McDonell, Niklas Muenighoff, Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron, Lintang Sutawika, Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. The language model evaluation harness, 07 2024. URL <https://zenodo.org/records/12608602>.
- Alex Graves, Greg Wayne, and Ivo Danihelka. Neural turing machines, 2014. URL <https://arxiv.org/abs/1410.5401>.
- Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. In *First conference on language modeling*, 2024.
- Albert Gu, Karan Goel, and Christopher Ré. Efficiently modeling long sequences with structured state spaces. *arXiv preprint arXiv:2111.00396*, 2021.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding, 2021. URL <https://arxiv.org/abs/2009.03300>.
- Angelos Katharopoulos, Apoorv Vyas, Nikolaos Pappas, and François Fleuret. Transformers are rnns: Fast autoregressive transformers with linear attention. In *International conference on machine learning*, pp. 5156–5165. PMLR, 2020.
- Lucas D Lingle. Transformer-vq: Linear-time transformers via vector quantization. *arXiv preprint arXiv:2309.16354*, 2023.
- Ilya Loshchilov and Frank Hutter. Sgdr: Stochastic gradient descent with warm restarts, 2017. URL <https://arxiv.org/abs/1608.03983>.
- Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization, 2019. URL <https://arxiv.org/abs/1711.05101>.
- Anton Lozhkov, Loubna Ben Allal, Leandro von Werra, and Thomas Wolf. Fineweb-edu: the finest collection of educational content, 2024. URL <https://huggingface.co/datasets/HuggingFaceFW/fineweb-edu>.
- Xuezhe Ma, Chunting Zhou, Xiang Kong, Junxian He, Liangke Gui, Graham Neubig, Jonathan May, and Luke Zettlemoyer. Mega: Moving average equipped gated attention. *arXiv preprint arXiv:2209.10655*, 2022.

- Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mixture models, 2016. URL <https://arxiv.org/abs/1609.07843>.
- Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. Can a suit of armor conduct electricity? a new dataset for open book question answering. *arXiv preprint arXiv:1809.02789*, 2018.
- Bo Peng, Eric Alcaide, Quentin Anthony, Alon Albalak, Samuel Arcadinho, Stella Biderman, Huanqi Cao, Xin Cheng, Michael Chung, Matteo Grella, et al. Rwkv: Reinventing rnns for the transformer era. *arXiv preprint arXiv:2305.13048*, 2023.
- Bo Peng, Ruichong Zhang, Daniel Goldstein, Eric Alcaide, Xingjian Du, Haowen Hou, Jiaju Lin, Jiaying Liu, Janna Lu, William Merrill, et al. Rwkv-7” goose” with expressive dynamic state evolution. *arXiv preprint arXiv:2503.14456*, 2025.
- Michael Poli, Stefano Massaroli, Eric Nguyen, Daniel Y Fu, Tri Dao, Stephen Baccus, Yoshua Bengio, Stefano Ermon, and Christopher Ré. Hyena hierarchy: Towards larger convolutional language models. In *International Conference on Machine Learning*, pp. 28043–28078. PMLR, 2023.
- Zhen Qin, Songlin Yang, Weixuan Sun, Xuyang Shen, Dong Li, Weigao Sun, and Yiran Zhong. Hgrn2: Gated linear rnns with state expansion, 2024. URL <https://arxiv.org/abs/2404.07904>.
- Jack W Rae, Anna Potapenko, Siddhant M Jayakumar, and Timothy P Lillicrap. Compressive transformers for long-range sequence modelling. *arXiv preprint arXiv:1911.05507*, 2019.
- Melissa Roemmele, Cosmin Adrian Bejan, and Andrew S Gordon. Choice of plausible alternatives: An evaluation of commonsense causal reasoning. In *2011 AAAI Spring Symposium Series*, 2011. URL <https://people.ict.usc.edu/~gordon/publications/AAAI-SPRING11A.PDF>.
- Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Winogrande: An adversarial winograd schema challenge at scale, 2019. URL <https://arxiv.org/abs/1907.10641>.
- Imanol Schlag, Kazuki Irie, and Jürgen Schmidhuber. Linear transformers are secretly fast weight programmers. In Marina Meila and Tong Zhang (eds.), *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pp. 9355–9366. PMLR, 18–24 Jul 2021a. URL <https://proceedings.mlr.press/v139/schlag21a.html>.
- Imanol Schlag, Kazuki Irie, and Jürgen Schmidhuber. Linear transformers are secretly fast weight programmers. In *International conference on machine learning*, pp. 9355–9366. PMLR, 2021b.
- Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan, Wen Bo, and Yunfeng Liu. Roformer: Enhanced transformer with rotary position embedding. *Neurocomputing*, 568:127063, 2024.
- Yu Sun, Xinhao Li, Karan Dalal, Jiarui Xu, Arjun Vikram, Genghan Zhang, Yann Dubois, Xinlei Chen, Xiaolong Wang, Sanmi Koyejo, et al. Learning to (learn at test time): Rnns with expressive hidden states. *arXiv preprint arXiv:2407.04620*, 2024.
- Yutao Sun, Li Dong, Shaohan Huang, Shuming Ma, Yuqing Xia, Jilong Xue, Jianyong Wang, and Furu Wei. Retentive network: A successor to transformer for large language models. *arXiv preprint arXiv:2307.08621*, 2023.
- Kimi Team, Yu Zhang, Zongyu Lin, Xingcheng Yao, Jiayi Hu, Fanqing Meng, Chengyin Liu, Xin Men, Songlin Yang, Zhiyuan Li, Wentao Li, Enzhe Lu, Weizhou Liu, Yanru Chen, Weixin Xu, Longhui Yu, Yejie Wang, Yu Fan, Longguang Zhong, Enming Yuan, Dehao Zhang, Yizhi Zhang, T. Y. Liu, Haiming Wang, Shengjun Fang, Weiran He, Shaowei Liu, Yiwei Li, Jianlin Su, Jiezhong Qiu, Bo Pang, Junjie Yan, Zhejun Jiang, Weixiao Huang, Bohong Yin, Jiacheng You, Chu Wei, Zhengtao Wang, Chao Hong, Yutian Chen, Guanduo Chen, Yucheng Wang, Huabin Zheng, Feng Wang, Yibo Liu, Mengnan Dong, Zheng Zhang, Siyuan Pan, Wenhao Wu, Yuhao Wu, Longyu

- Guan, Jiawen Tao, Guohong Fu, Xinran Xu, Yuzhi Wang, Guokun Lai, Yuxin Wu, Xinyu Zhou, Zhilin Yang, and Yulun Du. Kimi linear: An expressive, efficient attention architecture, 2025. URL <https://arxiv.org/abs/2510.26692>.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation language models, 2023. URL <https://arxiv.org/abs/2302.13971>.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- Sinong Wang, Belinda Z Li, Madian Khabsa, Han Fang, and Hao Ma. Linformer: Self-attention with linear complexity. *arXiv preprint arXiv:2006.04768*, 2020.
- Johannes Welbl, Nelson F. Liu, and Matt Gardner. Crowdsourcing multiple choice science questions, 2017. URL <https://arxiv.org/abs/1707.06209>.
- Songlin Yang and Yu Zhang. Fla: A triton-based library for hardware-efficient implementations of linear attention mechanism, January 2024. URL <https://github.com/fla-org/flash-linear-attention>.
- Songlin Yang, Jan Kautz, and Ali Hatamizadeh. Gated delta networks: Improving mamba2 with delta rule. *arXiv preprint arXiv:2412.06464*, 2024.
- Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a machine really finish your sentence?, 2019. URL <https://arxiv.org/abs/1905.07830>.
- Hailin Zhang, Xiaodong Ji, Yilin Chen, Fangcheng Fu, Xupeng Miao, Xiaonan Nie, Weipeng Chen, and Bin Cui. Pqcache: Product quantization-based kvcache for long context llm inference. *Proceedings of the ACM on Management of Data*, 3(3):1–30, 2025.
- Yu Zhang and Songlin Yang. Flame: Flash language modeling made easy, January 2025. URL <https://github.com/fla-org/flame>.

## A LANGUAGE MODELING EXPERIMENT

Table 2 details the experimental setup. For RAM-Net, we consistently use a Top-8 setting for both training and evaluation. Regarding positional encoding, we apply CAPE to all heads in the initial 4 layers and half of the heads in the subsequent 4 layers, while the remaining heads operate without positional embeddings. The complete training configuration is summarized in Table 3.

Model	Width hidden size	Layers number of stack	Heads attn. heads	Active State per token	Total State per token	Other Configure feature size
Transformer++	1024	24	16	50.3M	50.3M	$d_k = d_v = 64$
GLA	1024	24	4	3.1M	3.1M	$d_k = 128, d_v = 256$
HGRN2	1024	24	8	3.1M	3.1M	$d_f = d_i = 128$
GDN	1024	21	3	8.5M	8.5M	$d_h = 256, d_v = 512, k_{conv} = 4$
Mamba2	1024	48	32	12.9M	12.9M	$d_h = 64, d_{ssm} = 128, k_{conv} = 4$
RAM-Net (Top-8)	1024	27	16	0.4M	28.8M	$d_v = 64, U = 5, d_p = 4$

Table 2: **Comparison of architectural configurations and state overheads.** *Active State* denotes the effective state size utilized during per-token computation, whereas *Total State* represents the full state size. Here, Transformer++ is evaluated at sequence length 1024.

Category	Parameter	Value
<b>Data</b>	Dataset	FineWeb-Edu (Lozhkov et al., 2024)
	Training budget	10B tokens
<b>Tokenizer</b>	Tokenizer	Llama 2 (Touvron et al., 2023)
	Vocabulary size	32k
	Context window	4,096
<b>Optimization</b>	Optimizer	AdamW (Loshchilov & Hutter, 2019)
	Weight decay	0.1
	Gradient clipping	1.0
<b>LR Schedule</b>	Schedule method	Cosine decay (Loshchilov & Hutter, 2017)
	Linear Warm-up	500M tokens
	Peak LR	$1.0 \times 10^{-3}$
	Min LR	$1.0 \times 10^{-4}$
<b>Batch</b>	Batch per device	65.5K tokens
	Global batch size	0.5M tokens
	Device number	8

Table 3: Training configuration summary.

## B MEMORY ACCESS PATTERN

To investigate the memory interaction dynamics of RAM-Net, we visualize the memory access traces in Fig. 5. This figure depicts the slot activation patterns of four representative heads selected from different layers, where read and write events exhibit distinct diagonal shifts attributable to positional encoding offsets. Specifically, panels (a) and (b) correspond to heads utilizing CAPE. In contrast, panels (c) and (d) operate without CAPE, displaying position-agnostic patterns driven purely by semantic content. Collectively, these traces reveal a hybrid mechanism of fixed and dynamic accesses. The observed spatiotemporal locality in these patterns highlights significant potential for future system-level optimizations, such as caching strategies.

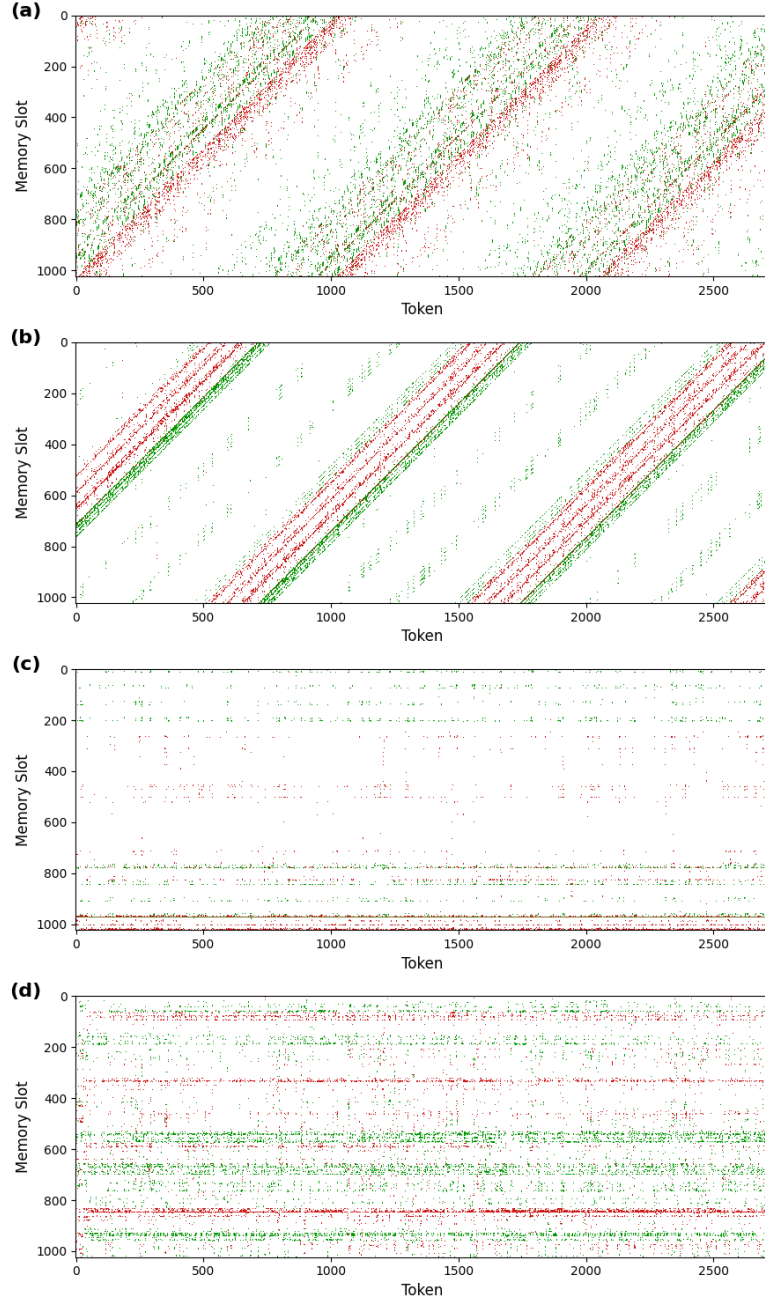


Figure 5: Visualization of memory access traces: read (green) and write (red) events across memory slots over time (tokens).