# Log-Linear Attention

**Han Guo**[1]* **Songlin Yang**[1]* **Tarushii Goel**[1] **Eric P. Xing**[3] **Tri Dao**[2] **Yoon Kim**[1]

[1]Massachusetts Institute of Technology [2]Princeton University, Together AI
[3]Carnegie Mellon University, Mohamed bin Zayed University of AI, GenBio AI

hanguo@mit.edu

## Abstract

The attention mechanism in Transformers is an important primitive for accurate and scalable sequence modeling. Its quadratic-compute and linear-memory complexity however remain significant bottlenecks. Linear attention and state-space models enable linear-time, constant-memory sequence modeling and can moreover be trained efficiently through matmul-rich parallelization across sequence length. However, at their core these models are still RNNs, and thus their use of a fixed-size hidden state to model the context is a fundamental limitation. This paper develops log-linear attention, an attention mechanism that balances linear attention's efficiency and the expressiveness of softmax attention. Log-linear attention replaces the fixed-size hidden state with a logarithmically growing set of hidden states. We show that with a particular growth function, log-linear attention admits a similarly matmul-rich parallel form whose compute cost is log-linear in sequence length. Log-linear attention is a general framework and can be applied on top of existing linear attention variants. As case studies, we instantiate log-linear variants of two recent architectures—Mamba-2 and Gated DeltaNet—and find they perform well compared to their linear-time variants.[2]

## 1 Introduction

The attention layer [4] is a core building block of modern deep learning architectures, most notably in the Transformer architecture [61]. For training, attention can be parallelized across sequence length through reformulating the computation as a series of matrix-matrix multiplications (matmuls), which can enable efficient training on modern accelerators such as GPUs and TPUs. However, the compute cost of attention grows quadratically and its memory cost grows linearly with respect to sequence length; despite the wallclock efficiency improvements obtained from hardware-optimized attention implementations [15, 11, 55, 34, 32], this quadratic-compute linear-memory cost is a fundamental limitation in enabling new applications and serves as a significant bottleneck in existing ones.

Linear attention [28] replaces the softmax kernel with a simple linear kernel (i.e., dot product) to derive the "attention" scores. The use of a linear kernel makes it possible to reformulate linear attention as a linear RNN with matrix-valued hidden states, and thus linear attention enables linear-time, constant-memory sequence modeling.[3] For training, linear attention can be parallelized across sequence length via a chunking mechanism where a sequence is split up into chunks and the computations across chunks are performed in parallel [23, 58, 67, 12]. The complexity of this chunkwise parallel algorithm is subquadratic in sequence length but still rich in matmuls,[4] leading to hardware-efficient implementations [64, 49, 6] that obtain practical wallclock improvements

---

*Equal contribution.
[2]Code available at https://github.com/HanGuo97/log-linear-attention.
[3]Thus there are three senses in which linear attention is *linear*: the use of a linear kernel, its reformulation as a linear RNN where the hidden state is a linear function of the previous state, and its linear-time complexity.
[4]Unlike parallel scan [8] which can also parallelize linear attention across sequence length but consists mostly of elementwise operations instead of matmuls.

footer_navigationPreprint. Under review.

| Model | A | M (Data Dependent?) | Training Algorithm / Time | Decoding Time and Space | |
|---|---|---|---|---|---|
| Attention | $\sigma(\mathbf{QK}^\top)$ | Mask (✗) | FlashAttention $\mathcal{O}(T^2)$ | $\mathcal{O}(T)$ | $\mathcal{O}(T)$ |
| Linear Attention [28] | $\mathbf{QK}^\top$ | Mask (✗) | Chunk-recurrent $\mathcal{O}(T)$ | $\mathcal{O}(1)$ | $\mathcal{O}(1)$ |
| RetNet [58] | $\mathbf{QK}^\top$ | Semiseparable (✗) | Chunk-recurrent $\mathcal{O}(T)$ | $\mathcal{O}(1)$ | $\mathcal{O}(1)$ |
| Mamba-2 [12] | $\mathbf{QK}^\top$ | Semiseparable (✓) | Chunk-recurrent $\mathcal{O}(T)$ | $\mathcal{O}(1)$ | $\mathcal{O}(1)$ |
| Multi-Hyena [36] | $\mathbf{QK}^\top$ | Toeplitz (✗) | FFT $\mathcal{O}(T \log T)$ | $\mathcal{O}(\log^2 T)$ | $\mathcal{O}(T)$ |
| DeltaNet [53, 68] | $\mathcal{T}_\mathbf{K}(\mathbf{QK}^\top)$ | Mask (✗) | Chunk-recurrent $\mathcal{O}(T)$ | $\mathcal{O}(1)$ | $\mathcal{O}(1)$ |
| Gated DeltaNet [66] | $\mathcal{T}_\mathbf{K}(\mathbf{QK}^\top)$ | Semiseparable (✓) | Chunk-recurrent $\mathcal{O}(T)$ | $\mathcal{O}(1)$ | $\mathcal{O}(1)$ |
| Log-Linear Mamba-2 | $\mathbf{QK}^\top$ | Hierarchical (✓) | Chunk-scan $\mathcal{O}(T \log T)$ | $\mathcal{O}(\log T)$ | $\mathcal{O}(\log T)$ |
| Log-Linear Gated DeltaNet | $\mathcal{T}_\mathbf{K}(\mathbf{QK}^\top)$ | Hierarchical (✓) | Chunk-scan $\mathcal{O}(T \log T)$ | $\mathcal{O}(\log T)$ | $\mathcal{O}(\log T)$ |

**Table 1:** Summary of efficient attention mechanisms under the unified formulation: $\mathbf{P} = \mathbf{A} \odot \mathbf{M}, \mathbf{O} = \mathbf{PV}$. $\mathbf{M}$ is a lower-triangle (causal) matrix. We use symbol $\mathcal{T}_\mathbf{K}(\mathbf{A}) = (\mathbf{A} \odot \mathbf{L})(\mathbf{I} + \mathbf{KK}^\top \odot (\mathbf{I} - \mathbf{L}))^{-1}$ for notational brevity, where $\mathbf{L}$ is a lower-triangular matrix of 1s. Here decoding time is the time per step, and decoding space refers to the overall memory complexity during generation.

over optimized implementations of softmax attention. While early versions of linear attention generally underperformed softmax attention [27, 44, 35, 47, 58], modern variants with data-dependent multiplicative gates [67, 50, 42]—which include state-space models (SSMs) such as Mamba [19, 12]—and delta-rule-based structured transition matrices [53, 67, 66, 18, 57] have led to significant improvements. However, despite these improvements linear attention's use of a fixed-sized hidden state is a fundamental limitations when it comes to certain capabilities such as associative recall over a given context [2]. And empirically, although many recent linear RNNs claim to match or outperform softmax attention, these results are often based on training and evaluation in short-context settings. In practice, their performance degrades when exposed to longer contexts [59, 33].

This paper develops *log-linear attention* as a middle ground between linear attention and full softmax attention. Instead of using a single hidden state matrix to represent the history (as in linear attention/SSMs), log-linear attention maintains a *growing set* of hidden states where the set size grows logarithmically with respect to sequence length. With a particular choice of the growth function, we show that log-linear attention admits a matmul-rich "parallel form" for training which involves replacing the lower-triangular causal mask in ordinary linear attention with a data-dependent *hierarchical matrix*, which enables subquadratic training; in particular we show that the compute cost of log-linear attention is log-linear in sequence length (hence the name), while its memory cost is logarithmic. Log-linear attention is a general framework for sequence modeling and can be used to generalize existing linear attention models. As case studies, we use the framework on two popular recent models, Mamba-2 [12] and Gated DeltaNet [66] to derive log-linear variants of both models, and find that these variants perform well compared to their original linear variants.

## 2   Background: A Structured Matrix View of Efficient Attention Variants

Given an input sequence of length $T$ and the corresponding key, query, value matrices $\mathbf{K}, \mathbf{Q}, \mathbf{V} \in \mathbb{R}^{T \times d}$, softmax attention obtains the output $\mathbf{O} \in \mathbb{R}^{T \times d}$ for all time steps via $\mathbf{O} = \text{softmax}(\mathbf{QK}^\top \odot \mathbf{M})\mathbf{V}$, where $\mathbf{M} \in \{-\infty, 0\}^{T \times T}$ is a causal masking matrix. This incurs $\mathcal{O}(T^2)$ compute and $\mathcal{O}(T)$ memory, which makes it costly to apply to long sequences. As a response, there has been much recent work on efficient alternatives with sub-quadratic compute and sub-linear memory, including linear attention, state-space models, and long convolution models. Despite their differences, many of these approaches can be captured by the following equation:

$$\mathbf{P} = \mathbf{A} \odot \mathbf{M}, \quad \mathbf{O} = \mathbf{PV}, \tag{1}$$

where $\mathbf{A} \in \mathbb{R}^{T \times T}$ is an attention-like matrix (e.g., $\mathbf{QK}^\top$ in the case of ordinary linear attention) and $\mathbf{M} \in \mathbb{R}^{T \times T}$ is a lower-triangular causal masking matrix (e.g., $\mathbf{M} \in \{0, 1\}^{T \times T}$ for linear attention). By separating out the interaction terms $\mathbf{A}$ and the (potentially data-dependent) masking matrix $\mathbf{M}$ (which typically models the "decay factor" between two positions), this abstraction reveals commonalities across a broad class of models, as shown in Table 1. Moreover, different structures imposed on $\mathbf{M}$ can lead to efficient training and inference algorithms. We now describe key models that fit within this framework.

**Linear attention.**   Linear attention [28] simply removes the softmax operation, resulting in the following parallel form[5]

$$\mathbf{O} = (\mathbf{QK}^\top \odot \mathbf{M})\mathbf{V}, \quad \mathbf{M}_{ij} = \mathbf{1}\{i \le j\}.$$

---

[5]Here we work linear attention without any feature maps or normalization, since most recent works have found them to be unnecessary (although see [25, 9, 2]).

Linear attention can be reparameterized into the following "recurrent form" for inference,

$$\mathbf{S}_t = \mathbf{S}_{t-1} + \boldsymbol{v}_t \boldsymbol{k}_t^\top, \quad \boldsymbol{o}_t = \mathbf{S}_t \boldsymbol{q}_t,$$

which enables linear-time constant-memory sequence modeling.

**Linear attention with (data-dependent) gates.** Vanilla linear attention lacks a forgetting mechanism, which has been shown to be crucial in ordinary RNNs. One way to incorporate such a mechanism is through a scalar gate $\alpha_t \in (0, 1)$, which results in recurrence $\mathbf{S}_t = \alpha_t \mathbf{S}_{t-1} + \boldsymbol{v}_t \boldsymbol{k}_t^\top$. This has the following corresponding parallel form:

$$\mathbf{O} = (\mathbf{Q}\mathbf{K}^\top \odot \mathbf{M})\mathbf{V}, \quad \mathbf{M}_{ij} = \prod_{k=j+1}^{i} \alpha_k. \tag{2}$$

Originally introduced by Peng et al. [44], gated linear attention has enjoyed a resurgence in recent years [50, 42, 65, 29] and are an instance of time-varying SSMs [19, 12]. Well-known models in this family include RetNet [58], which uses a data-*in*dependent gate $\alpha_t = \alpha$, and Mamba-2 [12], which uses the above data-dependent gate. Dao and Gu [12] show that with a scalar gating factor, $\mathbf{M}$ has a 1-semiseparable structure where every submatrix in the lower triangular portion has rank at most 1. Some works such as Mamba [19] and GLA [67] work with data-dependent gating *matrices*, i.e., $\mathbf{S}_t = \mathbf{G}_t \odot \mathbf{S}_{t-1} + \boldsymbol{v}_t \boldsymbol{k}_t^\top$ where $\mathbf{G}_t \in (0, 1)^{d \times d}$. When $\mathbf{G}_t$ has rank-one structure (e.g., as in GLA and GateLoop [29]), it is still possible to have an "efficient" representation via $\mathbf{O} = \left( \left( \left( (\mathbf{Q} \odot \mathbf{B}) (\mathbf{K}/\mathbf{B})^\top \right) \right) \odot \mathbf{M} \right) (\mathbf{V}/\mathbf{D}) \right) \odot \mathbf{D}$ where $\mathbf{B}, \mathbf{D} \in \mathbb{R}^{L \times d}$ capture the per-step low-rank factorizations of $\mathbf{G}_t$, and $\mathbf{M} \in \{0, 1\}^{T \times T}$ is now a simple causal masking matrix (see Yang et al. [67, §C] for the derivation). This type of representation is (to the best of our knowledge) not possible with the full rank $\mathbf{G}_t$ used by Mamba.

**Linear attention with the delta rule.** DeltaNet [53] is a type of linear attention layer which updates the hidden state via the delta rule [62],[6] where the recurrent form is given by[7]

$$\mathbf{S}_t = \mathbf{S}_{t-1} \left( \mathbf{I} - \boldsymbol{k}_t \boldsymbol{k}_t^\top \right) + \boldsymbol{v}_t \boldsymbol{k}_t^\top, \quad \boldsymbol{o}_t = \mathbf{S}_t \boldsymbol{q}_t.$$

While the original work used a purely recurrent form, Yang et al. [68] recently show that it is possible to parallelize DeltaNet across sequence length through leveraging a compact representation of Householder matrices [7, 24], resulting in the following parallel form (cf. [68, §3.2]):

$$\mathbf{O} = \left( \underbrace{\left( \mathbf{Q}\mathbf{K}^\top \odot \mathbf{L} \right) \left( \mathbf{I} + \mathbf{K}\mathbf{K}^\top \odot (\mathbf{L} - \mathbf{I}) \right)^{-1}}_{\mathbf{A}} \odot \mathbf{M} \right) \mathbf{V}$$

where $\mathbf{L}$ and $\mathbf{M}$ are lower-triangular matrices consisting of 1s. Since $\mathbf{A}$ itself is already lower-triangular, the causal masking matrix $\mathbf{M}$ is not strictly necessary in the above. However, by changing $\mathbf{M}$ to have 1-semiseparable structure as in Mamba-2, we can recover Gated DeltaNet [66], whose recurrence is given by $\mathbf{S}_t = \alpha_t \mathbf{S}_{t-1}(\mathbf{I} - \boldsymbol{k}_t \boldsymbol{k}_t^\top) + \boldsymbol{v}_t \boldsymbol{k}_t^\top$. Linear attention with such data-dependent structured transition matrices has been shown to be theoretically more expressive than linear attention with multiplicative gates when it comes to certain types of *state-tracking* tasks [38, 18, 57, 43], which make these layers attractive targets to generalize via our log-linear attention framework.

**Long convolution models.** Long-convolution sequence models, where the convolution kernel size equals the sequence length, can also be cast into this framework. For example Toeplitz neural network [48] and MultiHyena [36] layers are given by $\mathbf{O} = (\mathbf{Q}\mathbf{K}^\top \odot \mathbf{T}_h) \mathbf{V}$, where $\mathbf{T}_h$ is a causal Toeplitz matrix generated by a long convolution kernel $\boldsymbol{h} \in \mathbb{R}^T$, i.e., $\mathbf{T}_h[i, j] = \boldsymbol{h}[i - j]$ for $i \geq j$ and 0 otherwise. Other long convolutional variants like H3 [17] and Hyena [45] also admit a precise attention-style formulation, which has already been shown in past work [45, 36]. While the decoding speed of long convolution models can be improved from $\mathcal{O}(T)$ to $\mathcal{O}(\log^2 T)$ per step [41], their memory cost remains linear, i.e., the same as in softmax attention. However, some long convolution models such as S4 [20] admit a reparameterization into a time-invariant SSM and thus enjoy constant-memory inference. There has also been efforts to distill long convolution models into RNNs [36, 46], but these inherit the memory bottleneck of RNNs.

---

[6]Linear attention with the delta rule is also an instance of a fast-weight programmer [54].

[7]The actual DeltaNet recurrence is given by $\mathbf{S}_t = \mathbf{S}_{t-1}(\mathbf{I} - \beta_t \boldsymbol{k}_t \boldsymbol{k}_t^\top) + \boldsymbol{k}_t \boldsymbol{v}_t^\top$ where $\beta_t$ is a data-dependent scalar value in either $(0, 1)$ or $(0, 2)$, but we set $\beta_t = 1$ here for notational brevity.

**Relationship between masking structure and efficient algorithms.** Using an unstructured $\mathbf{M}$ (e.g., a random lower-triangular matrix) degrades both compute and memory complexity to softmax attention-levels, despite the absence of softmax; i.e., the *structure* of $\mathbf{M}$ is essential for training/inference efficiency, not just the removal of softmax. In linear attention where $\mathbf{M}$ is a lower-triangular matrix of 1's, we can compute $\mathbf{O}$ chunkwise, leading to an $\mathcal{O}(T)$ algorithm.[8] This algorithm generalizes to the gated case where $\mathbf{M}$ has 1-semiseparable structure as shown in the state-space duality framework [12]. Long convolution models can use FFT to bring down the cost to $\mathcal{O}(T \log T)$.

## 3 Log-Linear Attention

The preceding section shows that the structure of the masking matrix $\mathbf{M}$ in $\mathbf{O} = (\mathbf{A} \odot \mathbf{M})\mathbf{V}$ plays a key role in determining compute and memory costs. Our *log-linear attention* mechanism places a particular structure on $\mathbf{M}$ that enables the compute cost to be log-linear in $T$ (i.e., $\mathcal{O}(T \log T)$) and the memory cost to be logarithmic (i.e., $\mathcal{O}(\log T)$). Log-linear attention only modifies the masking matrix $\mathbf{M}$ and therefore can be used to generalize linear attention models whose $\mathbf{A}$ matrix can have different structure. As case studies, we show how to derive log-linear variants of Mamba-2 and Gated DeltaNet based on our framework.



**Figure 1:** Standard linear attention (top) vs. log-linear attention (bottom). The input consists of query, key, and value vectors.

Log-linear attention employs a Fenwick tree–based scheme [16] to hierarchically partition the input into power-of-two-sized segments. Each position summarizes a range ending at that point, enabling queries to attend to a logarithmic number of hidden states that capture past context at multiple temporal scales. This structure naturally emphasizes recent tokens through finer segmentation and supports $\mathcal{O}(\log T)$ time and space complexity during decoding. For training, we show that this formulation corresponds to a structured $\mathbf{M}$ that yields a parallel algorithm with $\mathcal{O}(T \log T)$ time and $\mathcal{O}(T)$ space complexity.
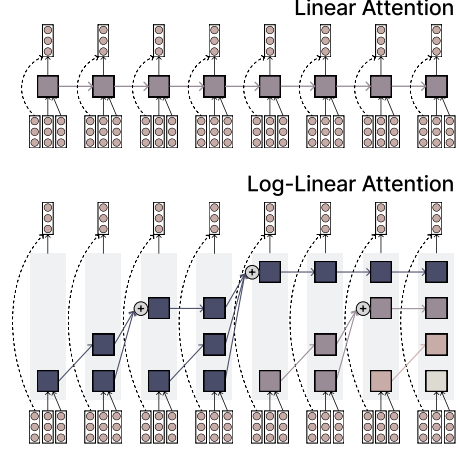
### 3.1 Fenwick Tree Partitioning for Linear Attention

We begin with the simplest form of linear attention and show how log-linear attention generalizes it by encoding distinct recurrent memories across different temporal segments.

For effective hierarchical segmentation, the method used to partition the prefix $[0, t)$ for a query $\boldsymbol{q}_t$ at step $t$ is critical. A straightforward approach assigns each token $s \in [t]$ to a level $\ell = \lfloor \log_2 s \rfloor$, based on its absolute position. However, in autoregressive decoding, this leads to overly coarse granularity for the most recent tokens—precisely the ones most crucial for accurate prediction. Intuitively, recent context should be modeled with higher resolution.



**Figure 2:** Fenwick tree bucket assignments.

To address this, we adopt a partitioning scheme based on the Fenwick tree structure [52, 16], which divides the prefix $[0, t)$ into up to $L = \lceil \log t \rceil + 1$ disjoint buckets. This decomposition is guided by the function $\text{lssb}(t) = \max\{\ell \in \mathbb{N} \mid 2^\ell \text{ divides } t\}$, which identifies the least significant set bit in the binary representation of $t$. Conceptually, the partitioning proceeds greedily, at each step subtracting the largest power of two that fits within the remaining segment of the prefix, asgiven below and shown in Fig. 2,
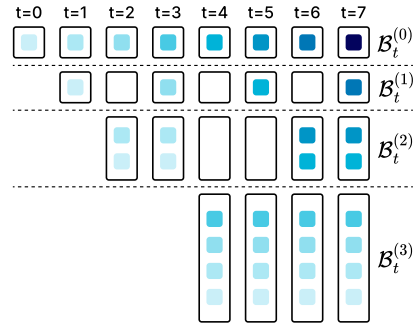
$$b_t^{(i)} = \begin{cases} t & \text{if } i = 0 \\ b_t^{(i-1)} - 2^{\text{lssb}\left(b_t^{(i-1)}\right)} & \text{otherwise} \end{cases} , \quad \mathcal{B}_t^{(\ell)} = \begin{cases} \{b_t^{(0)}\} & \text{if } \ell = 0 \\ \{b_t^{(i+1)}, \cdots, b_t^{(i)} - 1\} & \text{if } \ell = \text{lssb}\left(b_t^{(i)}\right) + 1 \\ \varnothing & \text{otherwise} \end{cases}$$

---

[8]This algorithm depends on the chunk size $C$, but since $C$ is a hyperparameter this is still linear in $T$.

Each bucket $\mathcal{B}_t^{(\ell)}$ has (at most) a power-of-two length: $|\mathcal{B}_t^{(\ell)}| = 2^{\ell-1}$ for $\ell \geq 1$, with a sentinel bucket of size $|\mathcal{B}_t^{(0)}| = 1$. Then, to obtain the output $\boldsymbol{o}_t$, log-linear attention computes the recurrent memory separately for each bucket, and weight the output by a data-dependent scalar $\lambda_t^{(\ell)} \geq 0$, which modulates the contribution of its corresponding bucket to the output. These weights are parameterized as functions of the input $\boldsymbol{x}_t$ via a linear projection, allowing the model to adaptively attend to different temporal scales. Concretely, the output is given by,

$$\boldsymbol{o}_t = \sum_{\ell=0}^{L-1} \lambda_t^{(\ell)} \boldsymbol{q}_t^\top \left( \sum_{s \in \mathcal{B}_t^{(\ell)}} \boldsymbol{v}_s \boldsymbol{k}_s^\top \right) = \sum_{\ell=0}^{L-1} \lambda_t^{(\ell)} \boldsymbol{q}_t^\top \mathbf{S}_t^{(\ell)}, \tag{3}$$

where $\mathbf{S}_t^{(\ell)} \in \mathbb{R}^{d \times d}$ is hidden state that summarizes all the information in level $\ell$. We observe that when all $\lambda_t^{(\ell)}$ are the same (or more generally when the $\lambda_t^{(\ell)}$ and $\lambda_t^{(\ell')}$ are linearly related across time) log-linear attention collapses to linear attention. Allowing distinct $\lambda_t^{(\ell)}$ is therefore essential for capturing multi-scale temporal structure.

**Parallel form.** While Eq. 3 is conceptually intuitive, it involves primarily matrix-vector products, which are inefficient on modern hardware optimized for matrix-matrix operations. To better leverage hardware acceleration and enable parallel computation across time steps, we reformulate the computation into a matmul-friendly form as in Sec. 2:

$$\mathbf{O} = \underbrace{\left(\mathbf{Q}\mathbf{K}^\top \odot \mathbf{M}^{\mathcal{H}}\right)}_{\mathbf{A}} \mathbf{V}, \quad \mathbf{M}_{ts}^{\mathcal{H}} = \begin{cases} \lambda_t^{\ell(t,s)} & \text{if } s \leq t \\ 0 & \text{otherwise} \end{cases} \tag{4}$$

Here, $\ell(t, s)$ denotes the level to which token $s$ belongs at time $t$ under the Fenwick tree partitioning. For brevity, we omit the explicit dependence on $(t, s)$ when the context is clear. Notably, the matrix $\mathbf{A}$ exhibits a structured low-rank pattern induced by the Fenwick tree partitioning, as shown below. In §3.2 we show how we can exploit this structure to derive a $\mathcal{O}(T \log T)$ parallel training algorithm.

$$\begin{bmatrix} \lambda_0^{(0)} \boldsymbol{q}_0^\top \boldsymbol{k}_0 & & & \\ \lambda_1^{(1)} \boldsymbol{q}_1^\top \boldsymbol{k}_0 & \lambda_1^{(0)} \boldsymbol{q}_1^\top \boldsymbol{k}_1 & & \\ \begin{bmatrix} \lambda_2^{(2)} \boldsymbol{q}_2 \\ \lambda_3^{(2)} \boldsymbol{q}_3 \end{bmatrix} \begin{bmatrix} \boldsymbol{k}_0 \\ \boldsymbol{k}_1 \end{bmatrix}^\top & \begin{matrix} \lambda_2^{(0)} \boldsymbol{q}_2^\top \boldsymbol{k}_2 \\ \lambda_3^{(1)} \boldsymbol{q}_3^\top \boldsymbol{k}_2 \quad \lambda_3^{(0)} \boldsymbol{q}_3^\top \boldsymbol{k}_3 \end{matrix} & & \\ \begin{bmatrix} \lambda_4^{(3)} \boldsymbol{q}_4 \\ \lambda_5^{(3)} \boldsymbol{q}_5 \\ \lambda_6^{(3)} \boldsymbol{q}_6 \\ \lambda_7^{(3)} \boldsymbol{q}_7 \end{bmatrix} \begin{bmatrix} \boldsymbol{k}_0 \\ \boldsymbol{k}_2 \\ \boldsymbol{k}_3 \\ \boldsymbol{k}_1 \end{bmatrix}^\top & & \begin{matrix} \lambda_4^{(0)} \boldsymbol{q}_4^\top \boldsymbol{k}_4 \\ \lambda_5^{(1)} \boldsymbol{q}_5^\top \boldsymbol{k}_4 \quad \lambda_5^{(0)} \boldsymbol{q}_5^\top \boldsymbol{k}_5 \\ \begin{bmatrix} \lambda_6^{(2)} \boldsymbol{q}_6 \\ \lambda_7^{(2)} \boldsymbol{q}_7 \end{bmatrix} \begin{bmatrix} \boldsymbol{k}_4 \\ \boldsymbol{k}_5 \end{bmatrix}^\top & \begin{matrix} \lambda_6^{(0)} \boldsymbol{q}_6^\top \boldsymbol{k}_6 \\ \lambda_7^{(1)} \boldsymbol{q}_7^\top \boldsymbol{k}_6 \quad \lambda_7^{(0)} \boldsymbol{q}_7^\top \boldsymbol{k}_7 \end{matrix} \end{matrix} \end{bmatrix}$$

**Memory-efficient decoding.** Incremental token-by-token decoding proceeds as follows. Recall that $\text{lssb}(t)$ determines the index of the least significant set bit in the binary representation of $t$. The new set of hidden states $\{\mathbf{S}_t^{(\ell)}\}_\ell$ are then given by the below equation.

$$\mathbf{S}_t^{(\ell)} = \begin{cases} \boldsymbol{v}_t \boldsymbol{k}_t^\top & \text{if } \ell = 0 \\ 0 & \text{if } 0 < \ell \leq \text{lssb}(t) \\ \sum_{\ell'=0}^{\ell-1} \mathbf{S}_{t-1}^{(\ell')} & \text{if } \ell = \text{lssb}(t)+1 \\ \mathbf{S}_{t-1}^{(\ell)} & \text{if } \ell > \text{lssb}(t)+1 \end{cases}$$

This recurrence reflects the core structure of the Fenwick-partitioned memory: at each timestep, the current memory term $\boldsymbol{v}_t \boldsymbol{k}_t^\top$ is inserted into the finest-resolution bucket ($\ell = 0$), while all buckets up to and including $\text{lssb}(t)$ are merged and promoted into a coarser-resolution bucket. When $t$ is a power of two, a new level is added. This maintains $\mathcal{O}(\log T)$ memory during inference. This decoding process follows the same principle as the online update and query mechanism in Fenwick trees.

**Remark.** The matrix $\mathbf{M}^{\mathcal{H}}$ (and $\mathbf{A}$) is a lower-triangular instance of a hierarchical ($\mathcal{H}$) matrix—specifically, of the HODLR (Hierarchically Off-Diagonal Low-Rank) type. When constructed using schemes like the Fenwick tree, it inherits the recursive partitioning and low-rank off-diagonal blocks that define $\mathcal{H}$ matrices. This establishes a direct connection between log-linear attention and hierarchical matrices: the attention operator corresponds to structured matrix multiplication with an $\mathcal{H}$ matrix. We refer to $\mathbf{M}^{\mathcal{H}}$ as a quasi-$\mathcal{H}$ matrix—a specialized class lying between general $\mathcal{H}$ and semiseparable matrices, designed to support $\mathcal{O}(\log T)$-space recurrence. See Section B.1 for details.
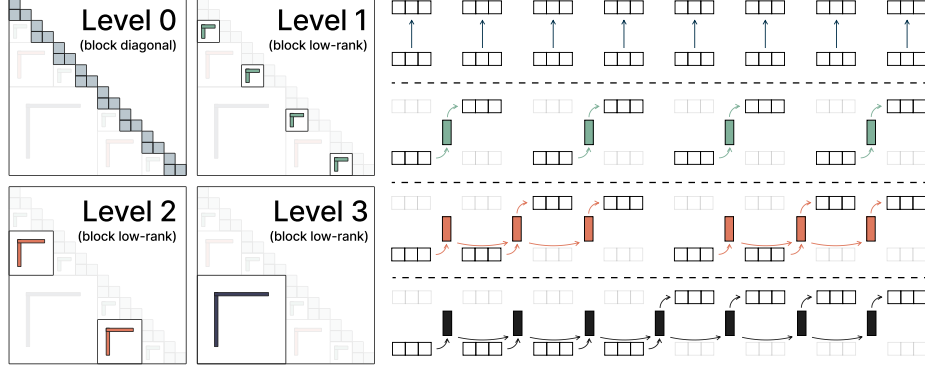
**Figure 3: Left**: Decomposition of the matrix $\mathbf{M}^{\mathcal{H}}$. **Right**: Chunkwise algorithm (Algorithm 1). Level 0 handles intra-chunk computations using a quadratic (in chunk size) algorithm, which is efficient due to small chunk sizes. Levels 1 and above perform inter-chunk computations by invoking existing inter-chunk primitives multiple times, with overall complexity logarithmic in the number of chunks.

## 3.2 Efficient Algorithm for Training

The chunkwise parallel algorithm for linear attention [58, 65, 12] splits a sequence into chunks of length $C$ and performs the computations for all chunks in parallel, while passing information across chunks when necessary. This offers a balance between the fully parallel and recurrent forms by reducing the computational cost of global attention while enabling greater sequence-level parallelism than strict recurrent computations. We show how primitives for efficient chunkwise computation of linear attention can be adapted to the log-linear case. First observe that the matrix $\mathbf{M}^{\mathcal{H}}$ exhibits a low-rank structure in its off-diagonal blocks, enabling a decomposition of the form:

$$\mathbf{M}^{\mathcal{H}} = \mathbf{D} + \sum_{\ell=1}^{L-1} \mathbf{M}^{(\ell)}, \quad \mathbf{M}^{(\ell)}_{ts} = \begin{cases} \lambda_t^{(\ell)}, & \text{if } s \in \mathcal{B}_t^{(\ell)}, \\ 0, & \text{otherwise.} \end{cases}$$

Here, $\mathbf{D}$ is a block-diagonal matrix with $\frac{T}{C}$ blocks $\{\mathbf{D}^{[k]}\}_{k=1}^{\frac{T}{C}}$, each encoding intra-chunk interactions. Each block $\mathbf{D}^{[i]} \in \mathbb{R}^{C \times C}$ is lower triangular, where $\mathbf{D}^{[i]}_{ts} = \lambda_{iC+t}^{(\ell)}$. $\mathbf{M}^{(\ell)}$ captures inter-chunk dependencies at level $\ell$ through a blockwise low-rank structure. See Fig. 3 (left) for an illustration.

Building on this structure, we develop a chunkwise log-linear attention algorithm (Algorithm 1). As shown in Fig. 3 (right), this chunkwise strategy adds a logarithmic overhead on top of linear attention. This algorithm processes the interactions in two stages.

**Intra-chunk computations** ($\ell=0$): For the block diagonal component $\mathbf{D}$, each block is treated as dense unstructured block, resulting in an ordinary $\mathcal{O}(C^2)$ matrix multiplication for the $\frac{T}{C}$ diagonal blocks, thus incurring $\mathcal{O}(TC)$ cost in total.

**Inter-chunk computations** ($\ell>0$): For the blocks corresponding to $\{\mathbf{M}^{(\ell)}\}_{\ell=1}^{L-1}$, the dependencies between chunks are handled via a sequence of linear attention passes. Owing to the hierarchical matrix structure and its decomposition (Eq. 3.2), each level reduces to a computation involving a sequentially semi-separable (SSS) matrix. When an efficient (linear-time) state-passing primitive is available—such as those used in Mamba-2 or Gated DeltaNet—inter-chunk computation requires only $\mathcal{O}(\log \frac{T}{C})$ invocations of this primitive. Each invocation costs $\mathcal{O}(T)$ in both time and memory,[9] and thus the total cost of these operations are $\mathcal{O}(T \log T)$.

This method extends the classical scan algorithm to the hierarchical domain, which we term a *chunkwise parallel scan*. Unlike token-level scans—often hindered by memory bandwidth limitations and high I/O overhead during training [65]—chunk scan restructures recurrent memory updates into parallel operations across chunks. Specifically, it performs $\mathcal{O}(\log T)$ independent scans, one per memory level, each implementable using standard parallel techniques such as the Blelloch scan [8]. Layer-specific weights (e.g., the $\lambda_t^{(\ell)}$ terms from $\mathbf{M}^{\mathcal{H}}$) are directly embedded in these scans, enabling efficient and scalable computation throughout the hierarchy.

---

[9]At level $\ell$, the matrix $\mathbf{M}^{(\ell)}$ contains $\frac{T}{2^{\ell-1}C}$ chunks, each of size $2^{\ell-1}C$. By skipping redundant operations, the total cost can be reduced by a constant factor of two.
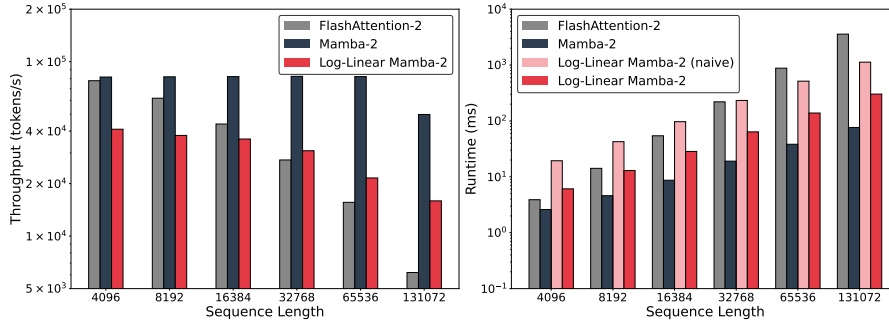
**Figure 4:** Training throughput (left; higher is better) and kernel runtime for a forward and backward pass (right; lower is better) across varying sequence lengths. **Log-Linear Mamba-2 (naive)** denotes repeated application of the existing Mamba-2 primitives, while **Log-Linear Mamba-2** uses a custom implementation with optimizations such as level fusion. The throughput drop at sequence length 131K is due to gradient checkpointing to reduce memory usage. Experiments were run on an H100 GPU with batch size 2, 48 heads, head dimension 64, state dimension 128, and chunk size 64. We use MVA for (Log-Linear) Mamba-2, and GQA for FlashAttention-2.

### 3.3 Log-Linear Generalizations of Mamba-2 and Gated DeltaNet

Having introduced log-linear attention as a hierarchical extension of linear attention, we now discuss how this idea can be applied to two concrete architectures: Mamba-2 [12] and Gated DeltaNet [66].

As discussed in Sec. 2, both models incorporate gating mechanisms, which induce a sequentially semi-separable (SSS) structure in the attention mask [12], which we denote as $\mathbf{M}^{\mathcal{S}}$ where $\mathbf{M}_{ij} = \prod_{k=j+1}^{i} \alpha_k$ (see Eq. 2). The main difference between the two lies in their respective parameterizations of the transition matrix $\mathbf{A}$. Our approach preserves the original form of $\mathbf{A}$ in each model while composing the attention mask with its log-linear variant $\mathbf{M} = \mathbf{M}^{\mathcal{S}} \odot \mathbf{M}^{\mathcal{H}}$.[10] We refer to the resulting models as *log-linear* Mamba-2 and *log-linear* Gated DeltaNet. Their parallel forms are given by,

$$\mathbf{O} = \left(\mathbf{Q}\mathbf{K}^T \odot \mathbf{M}^{\mathcal{S}} \odot \mathbf{M}^{\mathcal{H}}\right) \mathbf{V} \qquad \textit{Log-Linear Mamba-2}$$

$$\mathbf{O} = \left(\left(\mathbf{Q}\mathbf{K}^\top \odot \mathbf{L}\right)\left(\mathbf{I} + \mathbf{K}\mathbf{K}^\top \odot (\mathbf{L} - \mathbf{I})\right)^{-1} \odot \mathbf{M}^{\mathcal{S}} \odot \mathbf{M}^{\mathcal{H}}\right) \mathbf{V} \quad \textit{Log-Linear Gated DeltaNet}$$

This construction illustrates a general principle: any linear attention-style mechanism with structured memory and an efficient chunkwise-parallel primitive can be extended to a log-linear variant by composing its attention mask with a log-linear counterpart.

### 3.4 Implementation

We implemented the chunkwise parallel scan algorithm in `Triton` [60]. The custom kernel for log-linear Mamba-2 outperforms FlashAttention-2 [11] (forward + backward) at sequence lengths beyond 8K. In full training setups, throughput depends on model architecture. Notably, log-linear Mamba-2 (with MLP) surpasses Transformer throughput at 32K, despite additional layers like depthwise convolutions absent in the Transformer. See Sec. C for details.

### 3.5 Generalizations

**More expressive linear RNNs.** Consider general linear RNNs of the form $\mathbf{S}_t = \mathbf{S}_{t-1}\mathbf{C}_t + \boldsymbol{v}_t \boldsymbol{k}_t^\top$ where $\mathbf{C}_t$ is a data-dependent transition matrix. Without any structural assumptions on $\mathbf{C}_t$ (e.g., as in DeltaNet where $\mathbf{C}_t$ has identity-plus-rank-one structure) this general model does not admit a "nice" factorization of $\mathbf{A}$ into key/query matrices. This makes the general formulation difficult to work with in practice. However, it is still possible to derive log-linear versions of such expressive linear RNNs if we generalize $\mathbf{A}$ and $\mathbf{M}$ to higher-order tensors. We give this generalization in §A.
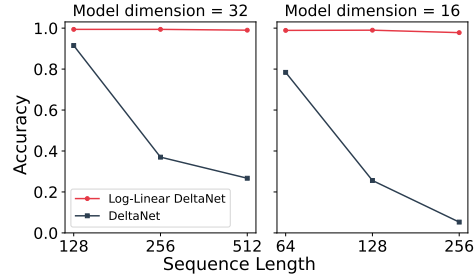
**Full use of levels.** As shown in Fig. 1, although there are $\mathcal{O}(\log T)$ states corresponding to different hierarchical levels, roughly half of them are zero in practice. This sparsity arises from the specific structure of HODLR matrices, which belong to a broader class of $\mathcal{H}$ matrices known as *weakly admissible* matrices [21]. Strongly admissible $\mathcal{H}$ matrices allow for finer-grained partitioning of the matrix, and their corresponding recurrent forms utilize all hierarchical levels. However, this comes at a significant computational cost, and we found this strongly admissible variant of log-linear attention to only marginally improve performance while incurring a significant slowdown. We thus adopt the weakly admissible structure throughout this work. See §B.4 for further discussion.

---

[10]More precisely, the elementwise product of an SSS matrix and an $\mathcal{H}$ matrix remains an $\mathcal{H}$ matrix. We separate them here for clarity.

# 4 Experiments

## 4.1 Synthetic Benchmark

We first experiment on multi-query associative recall (MQAR) [1], a standard diagnostic benchmark for evaluating the (in-context) recall capabilities of architectures. We train for 100 epochs on a dataset of 10K samples and sweep over learning rates. We only experiment with (non-gated) DeltaNet as this variant of linear attention performs best on MQAR. Our models use two layers, each with 1 head.



**Figure 5:** Experiments on MQAR. The number of key-value pairs in the input sequence is equal to the sequence length divided by 4.

**Results.** As shown in Fig. 5, we find that as the sequence length and number of key-value pairs increases, the performance of DeltaNet degrades significantly, while Log-Linear DeltaNet maintains high accuracy. Note that softmax attention obtains full accuracy on all settings.

## 4.2 Language Modeling

We perform academic-scale language modeling pretraining from scratch using 50B tokens on the Long-Data-Collections dataset,[11] using a sequence length of 16K. All models have 21 layers and use a hidden size of 1536. We use a Transformer with 16 attention heads and a RoPE base of 500K, a modified Mamba-2 with 48 heads and MLP layers, and a Gated DeltaNet with 6 heads. The Transformer, Mamba-2, and Gated DeltaNet models contain 693M, 802M, and 793M parameters, respectively. For the *log-linear* variants, we apply a linear layer on top of the hidden states to compute the per-head values $\lambda_t^{(\ell)}$. This adds less than $3\%$ additional parameters for Mamba-2 (825M) and less than $0.4\%$ for Gated DeltaNet (796M). Since Mamba-2 and Gated DeltaNet have more parameters than ordinary Transformers, we also include a (roughly) parameter-matched Transformer variant with 24 layers (778M parameters) for comparison. For our log-linaer variants, we use the default hyperparameters from the baselines (§D).

| Model | Wiki. ppl ↓ | LMB. ppl ↓ | LMB. acc ↑ | PIQA acc ↑ | Hella. acc_n ↑ | Wino. acc ↑ | ARC-e acc ↑ | ARC-c acc_n ↑ | Avg. |
|---|---|---|---|---|---|---|---|---|---|
| Transformer | 21.56 | 22.14 | 38.8 | 65.1 | 39.6 | 50.7 | 45.6 | 24.5 | 44.0 |
| w/ *24 Layers* | 21.13 | 21.17 | 39.3 | 66.6 | 40.4 | 53.3 | 47.8 | 26.4 | 45.6 |
| Mamba-2 | 22.44 | 24.14 | 36.2 | 66.8 | 41.2 | 51.6 | 46.0 | 27.1 | 44.8 |
| w/ *Log-Linear* | 22.11 | 21.86 | 37.0 | 66.6 | 41.1 | 51.7 | 45.5 | 27.4 | 44.9 |
| Gated DeltaNet | 21.73 | 19.71 | 39.3 | 65.8 | 40.9 | 52.2 | 47.1 | 24.6 | 45.0 |
| w/ *Log-Linear* | 21.44 | 18.08 | 40.5 | 66.1 | 41.4 | 53.9 | 46.9 | 24.9 | 45.6 |

**Table 2:** Performance comparison on language modeling and zero-shot commonsense reasoning.

**Standard benchmarks.** Following prior work [12, 66], we evaluate models on WikiText perplexity and several zero-shot commonsense reasoning benchmarks (Table 2). These are short-context tasks and are therefore largely insensitive to model state size. As such, we generally expect the log-linear variants to perform comparably to their linear counterparts. Log-Linear Mamba-2 improves upon its linear counterpart in perplexity and in half of the commonsense reasoning tasks. Log-Linear Gated DeltaNet shows stronger gains, outperforming its linear version in perplexity and in all but one reasoning benchmark. Notably, it also outperforms a layer-matched Transformer across all metrics and a parameter-matched Transformer on half of them.

**Per-position loss.** Following Lin et al. [33], we report the model's loss at each token position to evaluate its ability to handle long contexts (Fig. 6). If the loss steadily decreases as the token position increases, it suggests the model is effectively using the full context. However, if the loss levels off after a certain point, it indicates the model struggles to make use of information that is too far back
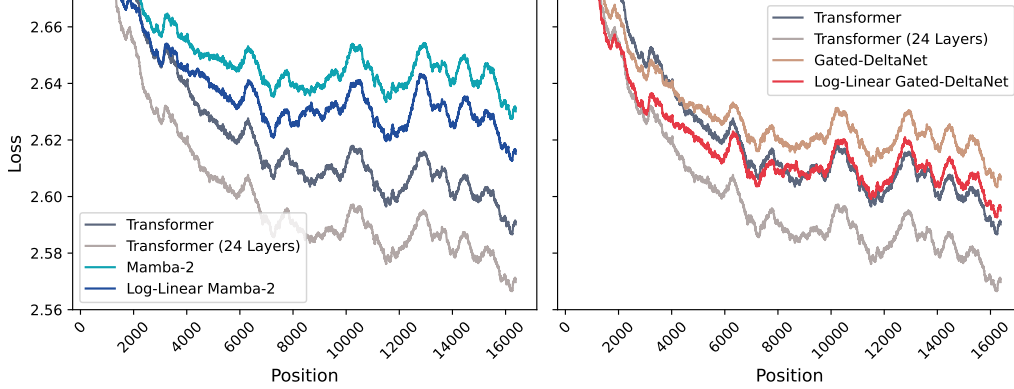
---

[11]https://huggingface.co/datasets/togethercomputer/Long-Data-Collections.

**Figure 6:** Per-position loss on Book3 samples (about 39M tokens) with running average of window size $501$.
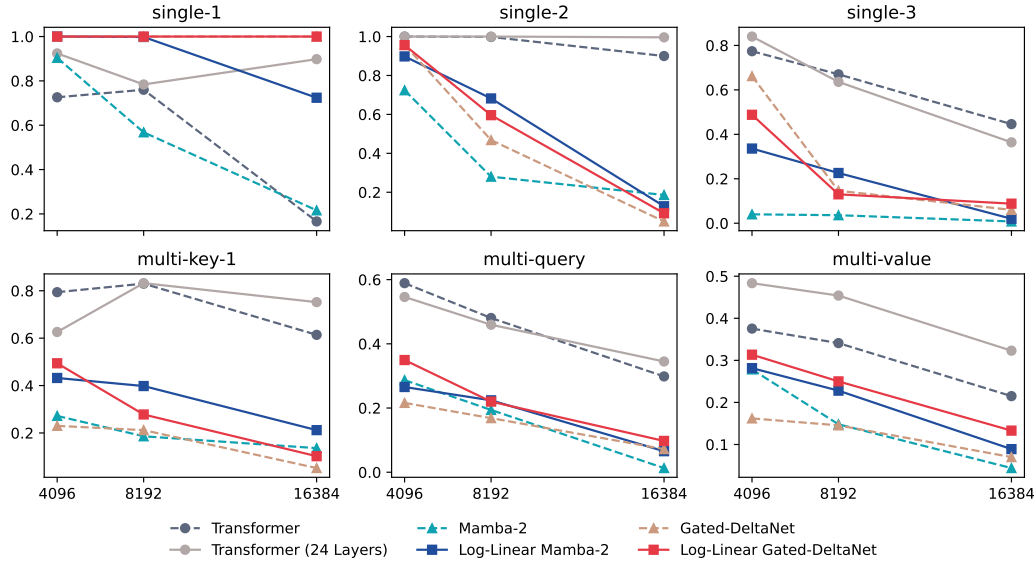


**Figure 7:** Needle-In-A-Haystack experiments. See Table 6 for details.

in the sequence. For this analysis, we use 39M tokens from Book-3.[12] To improve visualization, we apply a running average with a window size of 501. We observe that extending both Mamba-2 and Gated DeltaNet to their log-linear counterparts consistently reduces the (smoothed) loss across various positions, indicating improved long-range context utilization. Log-Linear Gated DeltaNet also closely tracks the performance of the layer-matched Transformer, although a performance gap remains when compared to the parameter-matched Transformer.

**Needle-In-A-Haystack.** We use the Needle-In-A-Haystack (NIAH, Fig. 7) benchmark from RULER [22], where the model must retrieve a value (the "needle") based on a key hidden in a long context (the "haystack"). In the simpler single-needle tasks, the log-linear variant of Mamba-2 outperformed its linear counterpart on 8 out of 9 metrics. Gated DeltaNet, which already achieved perfect accuracy in several cases, saw improvements in 3 metrics, with 3 remaining unchanged. For the more challenging multi-needle tasks, Log-Linear Mamba-2 again improved in 8 out of 9 metrics, while Log-Linear Gated DeltaNet achieved improvements across all metrics.

**In-Context Retrieval.** Following Arora et al. [2, 1], we evaluate models on real-world, recall-intensive tasks (Table 3). Since these benchmarks were originally designed for short sequences ($\leq$2K tokens), we report results at sequence lengths of 512, 1024, 2048, and (except NQ) 16K. We find that Log-Linear Mamba-2 yields improvements on roughly half of the tasks (SQuAD, TriviaQA, and NQ). In contrast, Log-Linear Gated DeltaNet shows more consistent gains, matching or outperforming Gated DeltaNet across all tasks except DROP.

---

[12]`victor-wu/book3`

**Table 3** (SWDE / SQuAD / FDA):

| Model | SWDE 512 | 1024 | 2048 | 16k | SQuAD 512 | 1024 | 2048 | 16k | FDA 512 | 1024 | 2048 | 16k |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Transformer | 47.3 | 44.6 | 45.2 | 45.4 | 34.0 | 34.5 | 34.5 | 34.5 | 72.2 | 70.8 | 72.9 | 72.2 |
| w/ *24 Layers* | 53.8 | 50.9 | 50.3 | 50.8 | 30.7 | 31.2 | 31.2 | 30.9 | 73.8 | 76.0 | 74.4 | 73.8 |
| Mamba-2 | 42.5 | 37.7 | 30.7 | 30.6 | 21.6 | 21.7 | 21.9 | 22.0 | 53.7 | 38.0 | 23.8 | 21.3 |
| w/ *Log-Linear* | 41.9 | 35.6 | 28.4 | 28.5 | 25.8 | 25.9 | 25.9 | 26.1 | 53.0 | 37.5 | 20.5 | 16.6 |
| Gated DeltaNet | 41.0 | 32.5 | 27.2 | 27.8 | 23.8 | 24.1 | 24.3 | 23.7 | 57.2 | 43.7 | 33.2 | 30.5 |
| w/ *Log-Linear* | 46.2 | 39.4 | 35.3 | 35.1 | 25.2 | 25.2 | 25.3 | 25.3 | 64.9 | 53.5 | 39.1 | 30.5 |

**Table 3** (TriviaQA / Drop / NQ):

| Model | TriviaQA 512 | 1024 | 2048 | 16k | Drop 512 | 1024 | 2048 | 16k | NQ 512 | 1024 | 2048 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Transformer | 48.5 | 49.6 | 48.5 | 48.5 | 22.8 | 22.8 | 22.5 | 22.3 | 24.5 | 24.3 | 24.6 |
| w/ *24 Layers* | 46.9 | 47.0 | 46.8 | 46.8 | 22.7 | 22.4 | 22.7 | 23.0 | 24.0 | 24.4 | 24.5 |
| Mamba-2 | 43.7 | 43.2 | 43.2 | 43.2 | 22.2 | 22.1 | 22.2 | 22.1 | 18.5 | 16.5 | 16.5 |
| w/ *Log-Linear* | 44.9 | 45.0 | 45.5 | 45.5 | 20.2 | 20.6 | 20.3 | 19.9 | 20.0 | 19.9 | 20.4 |
| Gated DeltaNet | 45.6 | 45.6 | 45.6 | 45.6 | 21.1 | 21.7 | 21.4 | 21.8 | 20.1 | 18.4 | 18.7 |
| w/ *Log-Linear* | 45.9 | 45.6 | 46.0 | 46.0 | 20.7 | 20.8 | 20.8 | 21.0 | 22.5 | 21.8 | 21.3 |

**Table 3:** Accuracy on retrieval tasks w/ input truncated to different lengths.

| Model | Single-Doc QA NQA | QQA | MFQ | Multi-Doc QA HQA | 2WM | Mus | Summarization GvR | QMS | MNs | Few-shot TRC | TQA | SSM | Code LCC | RBP |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Transformer | 11.7 | 9.7 | 20.8 | 22.4 | 29.8 | 6.7 | 13.1 | 9.4 | 3.2 | 27.5 | 28.0 | 16.2 | 23.7 | 29.8 |
| w/ *24 Layers* | 10.7 | 18.4 | 26.1 | 33.7 | 25.7 | 11.6 | 16.8 | 9.4 | 10.3 | 16.5 | 45.2 | 14.3 | 31.5 | 30.9 |
| Mamba-2 | 9.1 | 17.4 | 10.9 | 11.2 | 20.9 | 4.3 | 8.3 | 6.0 | 4.9 | 2.0 | 22.6 | 8.8 | 38.1 | 34.6 |
| w/ *Log-Linear* | 9.8 | 9.6 | 15.4 | 11.5 | 22.0 | 5.1 | 5.4 | 11.1 | 4.5 | 16.5 | 21.6 | 14.9 | 31.2 | 30.3 |
| Gated DeltaNet | 8.5 | 11.9 | 16.4 | 14.4 | 24.5 | 6.6 | 9.2 | 11.7 | 11.6 | 36.5 | 25.3 | 23.1 | 31.1 | 31.1 |
| w/ *Log-Linear* | 9.9 | 6.1 | 17.6 | 17.7 | 25.2 | 7.5 | 5.5 | 11.9 | 1.9 | 8.0 | 41.1 | 23.2 | 28.3 | 29.6 |

**Table 4:** Accuracy on LongBench tasks [5]: Narrative QA, QasperQA, MultiField QA, HotpotQA, 2WikiMultiQA, Musique, GovReport, QMSum, MultiNews, TREC, TriviaQA, SamSum, LCC, and RepoBench-P.

**Long context understanding.** Finally, we evaluated the models' performance on LongBench [5] (Table 4). We observe that both *Log-Linear* Mamba-2 and Gated DeltaNet outperforms the baseline Mamba-2 Gated DeltaNet in 8 out of 14 evaluation tasks.

## 5 Discussion and Limitations

While log-linear attention improves upon linear attention in many cases, there are still quite a few tasks where it did not improve upon the linear attention baselines. Due to compute resources we were unable to experiment with different parameterizations of the $\lambda$ terms (or hyperparameters in general), and it is possible that optimal parameterization of $\lambda$ could lead to improved results. We also still observe a significant performance gap compared to Transformers also persists across all benchmarks.

The engineering complexity of log-linear attention is higher. Inter-chunk computations conceptually resemble multiple applications of linear attention primitives, but intra-chunk operations require bespoke implementations. These intra-chunk mechanisms are a primary factor behind the speed differences. Additionally, the backward pass is more intricate, as it requires (manually) computing the gradients not only for the standard attention components but also for the additional $\lambda$ terms.

Finally, the use of Fenwick-tree partitioning (§3.1) introduces an inductive bias: recent tokens are allocated more fine-grained memory, while distant tokens are compressed more aggressively. This design reflects a natural assumption rooted in hierarchical matrix which posits that interactions between distant elements can be approximated in low-rank form. While intuitive and inspired by physical phenomena, this inductive bias may not be optimal for all applications. Future work could explore extensions that enable more flexible structures while preserving computational efficiency.

## 6 Related Work

Matrix multiplication serves as the computational backbone of modern deep learning. Contemporary architectures typically consist of a token mixing layer and a channel mixing layer, both of which

heavily depend on matrix multiplications. A growing body of research has investigated replacing dense matrices with *structured matrices*. For channel mixing, efforts include Butterfly matrices [13], Monarch matrices [14], and more recently, Block Tensor-Train matrices [51]. Token mixing has been exemplified by the family of linear attention models [28] and their various kernelizations [63]. Dao and Gu [12] generalize these approaches by extending low-rank structures to semi-separable matrices, enabling efficient recurrent inference and subsuming many recent recurrent models. Another line of work employs sparse attention patterns such as sliding-window attention (SWA), and several hybrid approaches have also emerged [40, 3, 39]. In this work, we introduce a hierarchical matrix formulation that supports state expansion while maintaining hardware-efficient training and inference.

Several prior efforts have focused on reducing the quadratic cost of attention to log-linear time complexity [30, 56, 10, 48, 17]. Reformer [30] employs locality-sensitive hashing (LSH) to efficiently cluster similar queries and keys. Multi-resolution attention [69] adopts a hierarchical approach, progressively refining attention scores from coarse to fine granularity, while Fast Multipole Attention [26] adapts the classical fast multipole method to efficiently model long-range interactions. In our work, we leverage the Fenwick tree data structure—a specialized binary indexed tree that enables efficient prefix sum calculations and updates in logarithmic time—to design an efficient attention layer during both training and decoding phases. While Zhu and Soricut [71] also employ hierarchical matrices for attention, their formulation is fully parallel and targeted at modest sequence lengths. In contrast, our approach adopts a chunkwise-parallel strategy with a custom Triton implementation optimized for long-sequence training. Derived from linear attention, our design imposes a structured $\mathcal{H}$-matrix constraint that ensures $\mathcal{O}(\log T)$ inference and $\mathcal{O}(T \log T)$ training complexity.

# 7 Conclusion

We introduced Log-Linear Attention, a general framework that extends a broad class of linear attention and state-space models to their log-linear counterparts—models with logarithmically growing state size. This framework offers both theoretical insights and practical benefits, linking structured matrix theory with hardware-efficient computation. As a case study, we applied this approach to two recent architectures: Mamba-2 and Gated DeltaNet.

## Acknowledgments

## References

[1] S. Arora, S. Eyuboglu, A. Timalsina, I. Johnson, M. Poli, J. Zou, A. Rudra, and C. Ré. Zoology: Measuring and improving recall in efficient language models. *arXiv preprint arXiv:2312.04927*, 2023.

[2] S. Arora, S. Eyuboglu, M. Zhang, A. Timalsina, S. Alberti, D. Zinsley, J. Zou, A. Rudra, and C. Ré. Simple linear attention language models balance the recall-throughput tradeoff. In *Proceedings of ICML*, 2024.

[3] S. Arora, S. Eyuboglu, M. Zhang, A. Timalsina, S. Alberti, D. Zinsley, J. Zou, A. Rudra, and C. Ré. Simple linear attention language models balance the recall-throughput tradeoff, 2025. URL https://arxiv.org/abs/2402.18668.

[4] D. Bahdanau, K. Cho, and Y. Bengio. Neural machine translation by jointly learning to align and translate. In *Proceedings of ICLR*, 2014.

[5] Y. Bai, X. Lv, J. Zhang, H. Lyu, J. Tang, Z. Huang, Z. Du, X. Liu, A. Zeng, L. Hou, et al. Longbench: A bilingual, multitask benchmark for long context understanding. *arXiv preprint arXiv:2308.14508*, 2023.

[6] M. Beck, K. Pöppel, P. Lippe, and S. Hochreiter. Tiled flash linear attention: More efficient linear rnn and xlstm kernels. *arXiv preprint arXiv:2503.14376*, 2025.

[7] C. H. Bischof and C. V. Loan. The WY representation for products of householder matrices. In *SIAM Conference on Parallel Processing for Scientific Computing*, 1985. URL `https://api.semanticscholar.org/CorpusID:36094006`.

[8] G. E. Blelloch. Prefix sums and their applications. 1990.

[9] J. Buckman, C. Gelada, and S. Zhang. Symmetric Power Transformers.

[10] H. J. Cunningham, G. Giannone, M. Zhang, and M. P. Deisenroth. Reparameterized multi-resolution convolutions for long sequence modelling. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL `https://openreview.net/forum?id=RwgNbIpCpk`.

[11] T. Dao. FlashAttention-2: Faster attention with better parallelism and work partitioning. In *Proceedings of ICLR*, 2024.

[12] T. Dao and A. Gu. Transformers are SSMs: Generalized models and efficient algorithms through structured state space duality. In *Proceedings of ICML*, 2024.

[13] T. Dao, A. Gu, M. Eichhorn, A. Rudra, and C. Ré. Learning fast algorithms for linear transforms using butterfly factorizations, 2020. URL `https://arxiv.org/abs/1903.05895`.

[14] T. Dao, B. Chen, N. S. Sohoni, A. D. Desai, M. Poli, J. Grogan, A. Liu, A. Rao, A. Rudra, and C. Ré. Monarch: Expressive structured matrices for efficient and accurate training. In K. Chaudhuri, S. Jegelka, L. Song, C. Szepesvári, G. Niu, and S. Sabato, editors, *International Conference on Machine Learning, ICML 2022, 17-23 July 2022, Baltimore, Maryland, USA*, volume 162 of *Proceedings of Machine Learning Research*, pages 4690–4721. PMLR, 2022. URL `https://proceedings.mlr.press/v162/dao22a.html`.

[15] T. Dao, D. Y. Fu, S. Ermon, A. Rudra, and C. Ré. FlashAttention: Fast and memory-efficient exact attention with IO-awareness. In *Proceedings of NeurIPS*, 2022.

[16] P. M. Fenwick. A new data structure for cumulative frequency tables. *Software: Practice and Experience*, 24, 1994. URL `https://api.semanticscholar.org/CorpusID:7519761`.

[17] D. Y. Fu, T. Dao, K. K. Saab, A. W. Thomas, A. Rudra, and C. Re. Hungry hungry hippos: Towards language modeling with state space models. In *The Eleventh International Conference on Learning Representations*, 2023. URL `https://openreview.net/forum?id=COZDyOWYGg`.

[18] R. Grazzi, J. Siems, J. K. Franke, A. Zela, F. Hutter, and M. Pontil. Unlocking state-tracking in linear RNNs through negative eigenvalues. In *Proceedings of ICLR*, 2025.

[19] A. Gu and T. Dao. Mamba: Linear-time sequence modeling with selective state spaces. In *Proceedings of CoLM*, 2024.

[20] A. Gu, K. Goel, and C. Ré. Efficiently modeling long sequences with structured state spaces. In *Proceedings of ICLR*, 2022.

[21] W. Hackbusch, B. N. Khoromskij, and R. Kriemann. Hierarchical matrices based on a weak admissibility criterion. *Computing*, 73:207–243, 2004.

[22] C.-P. Hsieh, S. Sun, S. Kriman, S. Acharya, D. Rekesh, F. Jia, Y. Zhang, and B. Ginsburg. Ruler: What's the real context size of your long-context language models? *arXiv preprint arXiv:2404.06654*, 2024.

[23] W. Hua, Z. Dai, H. Liu, and Q. Le. Transformer quality in linear time. In *International conference on machine learning*, pages 9099–9117. PMLR, 2022.

[24] T. Joffrain, T. M. Low, E. S. Quintana-Ortí, R. A. van de Geijn, and F. G. V. Zee. Accumulating householder transformations, revisited. *ACM Trans. Math. Softw.*, 32:169–179, 2006. URL `https://api.semanticscholar.org/CorpusID:15723171`.

[25] P. Kacham, V. Mirrokni, and P. Zhong. Polysketchformer: Fast transformers via sketching polynomial kernels. *arXiv preprint arXiv:2310.01655*, 2023.

[26] Y. Kang, G. Tran, and H. D. Sterck. Fast multipole attention: A divide-and-conquer attention mechanism for long sequences, 2024. URL `https://arxiv.org/abs/2310.11960`.

[27] J. Kasai, H. Peng, Y. Zhang, D. Yogatama, G. Ilharco, N. Pappas, Y. Mao, W. Chen, and N. A. Smith. Finetuning pretrained transformers into rnns. In *Proceedings of EMNLP*, 2021.

[28] A. Katharopoulos, A. Vyas, N. Pappas, and F. Fleuret. Transformers are rnns: Fast autoregressive transformers with linear attention. In *Proceedings of ICML*, 2020.

[29] T. Katsch. Gateloop: Fully data-controlled linear recurrence for sequence modeling. *arXiv preprint arXiv:2311.01927*, 2023.

[30] N. Kitaev, Ł. Kaiser, and A. Levskaya. Reformer: The efficient transformer. In *Proceedings of ICLR*, 2020.

[31] D. Kressner, S. Massei, and L. Robol. Low-rank updates and a divide-and-conquer method for linear matrix equations. *SIAM Journal on Scientific Computing*, 41(2):A848–A876, 2019.

[32] W. Kwon, Z. Li, S. Zhuang, Y. Sheng, L. Zheng, C. H. Yu, J. Gonzalez, H. Zhang, and I. Stoica. Efficient memory management for large language model serving with pagedattention. In *Proceedings of SOSP*, 2023.

[33] Z. Lin, E. Nikishin, X. He, and A. Courville. Forgetting transformer: Softmax attention with a forget gate. In *The Thirteenth International Conference on Learning Representations*, 2025. URL `https://openreview.net/forum?id=q2Lnyegkr8`.

[34] H. Liu, M. Zaharia, and P. Abbeel. Ring attention with blockwise transformers for near-infinite context. In *Proceedings of ICLR*, 2024.

[35] H. H. Mao. Fine-Tuning Pre-trained Transformers into Decaying Fast Weights. In *Proceedings of EMNLP*, pages 10236–10242.

[36] S. Massaroli, M. Poli, D. Y. Fu, H. Kumbong, R. N. Parnichkun, D. W. Romero, A. Timalsina, Q. McIntyre, B. Chen, A. Rudra, C. Zhang, C. Re, S. Ermon, and Y. Bengio. Laughing hyena distillery: Extracting compact recurrences from convolutions. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL `https://openreview.net/forum?id=OWELckerm6`.

[37] S. Massei, L. Robol, and D. Kressner. hm-toolbox: Matlab software for hodlr and hss matrices. *SIAM Journal on Scientific Computing*, 42(2):C43–C68, 2020.

[38] W. Merrill, J. Petty, and A. Sabharwal. The illusion of state in state-space models. In *Proceedings of ICML*, 2024.

[39] T. Munkhdalai, M. Faruqui, and S. Gopal. Leave no context behind: Efficient infinite context transformers with infini-attention, 2024. URL `https://arxiv.org/abs/2404.07143`.

[40] T. Nguyen, V. Suliafu, S. Osher, L. Chen, and B. Wang. Fmmformer: Efficient and flexible transformer via decomposed near-field and far-field attention. In *Proceedings of NeurIPS*, 2021.

[41] C.-A. Oncescu, S. Purandare, S. Idreos, and S. M. Kakade. Flash inference: Near linear time inference for long convolution sequence models and beyond. In *The Thirteenth International Conference on Learning Representations*, 2025. URL `https://openreview.net/forum?id=cZWCjan02B`.

[42] B. Peng, D. Goldstein, Q. Anthony, A. Albalak, E. Alcaide, S. Biderman, E. Cheah, T. Ferdinan, H. Hou, P. Kazienko, et al. Eagle and finch: Rwkv with matrix-valued states and dynamic recurrence. *arXiv preprint arXiv:2404.05892*, 3, 2024.

[43] B. Peng, R. Zhang, D. Goldstein, E. Alcaide, X. Du, H. Hou, J. Lin, J. Liu, J. Lu, W. Merrill, et al. Rwkv-7" goose" with expressive dynamic state evolution. *arXiv preprint arXiv:2503.14456*, 2025.

[44] H. Peng, N. Pappas, D. Yogatama, R. Schwartz, N. A. Smith, and L. Kong. In *Proceedings of ICLR*, 2021.

[45] M. Poli, S. Massaroli, E. Nguyen, D. Y. Fu, T. Dao, S. Baccus, Y. Bengio, S. Ermon, and C. Ré. Hyena hierarchy: Towards larger convolutional language models. In A. Krause, E. Brunskill, K. Cho, B. Engelhardt, S. Sabato, and J. Scarlett, editors, *International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA*, volume 202 of *Proceedings of Machine Learning Research*, pages 28043–28078. PMLR, 2023. URL `https://proceedings.mlr.press/v202/poli23a.html`.

[46] Z. Qin and Y. Zhong. Accelerating toeplitz neural network with constant-time inference complexity. In H. Bouamor, J. Pino, and K. Bali, editors, *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 12206–12215, Singapore, Dec. 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.750. URL `https://aclanthology.org/2023.emnlp-main.750/`.

[47] Z. Qin, W. Sun, H. Deng, D. Li, Y. Wei, B. Lv, J. Yan, L. Kong, and Y. Zhong. cosformer: Rethinking softmax in attention. In *Proceedings of ICLR*, 2022.

[48] Z. Qin, X. Han, W. Sun, B. He, D. Li, D. Li, Y. Dai, L. Kong, and Y. Zhong. Toeplitz neural network for sequence modeling. In *Proceedings of ICLR*, 2023.

[49] Z. Qin, W. Sun, D. Li, X. Shen, W. Sun, and Y. Zhong. Lightning attention-2: A free lunch for handling unlimited sequence lengths in large language models. *arXiv preprint arXiv:2401.04658*, 2024.

[50] Z. Qin, S. Yang, W. Sun, X. Shen, D. Li, W. Sun, and Y. Zhong. HGRN2: Gated Linear RNNs with State Expansion. In *Proceedings of CoLM*, 2024.

[51] S. Qiu, A. Potapczynski, M. Finzi, M. Goldblum, and A. G. Wilson. Compute better spent: Replacing dense layers with structured matrices. *ArXiv*, abs/2406.06248, 2024. URL `https://api.semanticscholar.org/CorpusID:270371652`.

[52] B. Y. Ryabko. A fast on-line adaptive code. *IEEE Trans. Inf. Theory*, 38:1400–1404, 1992. URL `https://api.semanticscholar.org/CorpusID:206392294`.

[53] I. Schlag, K. Irie, and J. Schmidhuber. Linear Transformers Are Secretly Fast Weight Programmers. In *Proceedings of ICML*, 2021.

[54] J. Schmidhuber. Learning to control fast-weight memories: An alternative to dynamic recurrent networks. *Neural Computation*, 4(1):131–139, 1992.

[55] J. Shah, G. Bikshandi, Y. Zhang, V. Thakkar, P. Ramani, and T. Dao. FlashAttention-3: Fast and accurate attention with asynchrony and low-precision. In *Proceedings of NeurIPS*, 2024.

[56] J. Shi, K. A. Wang, and E. B. Fox. Sequence modeling with multiresolution convolutional memory, 2023. URL `https://arxiv.org/abs/2305.01638`.

[57] J. Siems, T. Carstensen, A. Zela, F. Hutter, M. Pontil, and R. Grazzi. Deltaproduct: Improving state-tracking in linear rnns via householder products. *arXiv preprint arXiv:2502.10297*, 2025.

[58] Y. Sun, L. Dong, S. Huang, S. Ma, Y. Xia, J. Xue, J. Wang, and F. Wei. Retentive network: A successor to transformer for large language models. *arXiv preprint arXiv:2307.08621*, 2023.

[59] Y. Sun, X. Li, K. Dalal, J. Xu, A. Vikram, G. Zhang, Y. Dubois, X. Chen, X. Wang, S. Koyejo, T. Hashimoto, and C. Guestrin. Learning to (learn at test time): Rnns with expressive hidden states, 2025. URL `https://arxiv.org/abs/2407.04620`.

[60] P. Tillet, H.-T. Kung, and D. Cox. Triton: an intermediate language and compiler for tiled neural network computations. In *Proceedings of the 3rd ACM SIGPLAN International Workshop on Machine Learning and Programming Languages*, pages 10–19, 2019.

[61] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin. Attention is all you need. In *Proceedings of NeurIPS*, 2017.

[62] B. Widrow, M. E. Hoff, et al. Adaptive switching circuits. In *IRE WESCON convention record*, volume 4, pages 96–104. New York, 1960.

[63] Y. Xiong, Z. Zeng, R. Chakraborty, M. Tan, G. Fung, Y. Li, and V. Singh. Nyströmformer: A nyström-based algorithm for approximating self-attention. In *Proceedings of AAAI*, 2021.

[64] S. Yang and Y. Zhang. Fla: A triton-based library for hardware-efficient implementations of linear attention mechanism, Jan. 2024. URL `https://github.com/fla-org/flash-linear-attention`.

[65] S. Yang, B. Wang, Y. Shen, R. Panda, and Y. Kim. Gated linear attention transformers with hardware-efficient training. *arXiv preprint arXiv:2312.06635*, 2023.

[66] S. Yang, J. Kautz, and A. Hatamizadeh. Gated delta networks: Improving mamba2 with delta rule. In *Proceedings of ICLR*, 2024.

[67] S. Yang, B. Wang, Y. Shen, R. Panda, and Y. Kim. In *Proceedings of ICML*, 2024.

[68] S. Yang, B. Wang, Y. Zhang, Y. Shen, and Y. Kim. Parallelizing linear transformers with the delta rule over sequence length. In *Proceedings of NeurIPS*, 2024.

[69] Z. Zeng, S. Pal, J. Kline, G. M. Fung, and V. Singh. Multi resolution analysis (mra) for approximate self-attention, 2022. URL `https://arxiv.org/abs/2207.10284`.

[70] Y. Zhang and S. Yang. Flame: Flash language modeling made easy, Jan. 2025. URL `https://github.com/fla-org/flame`.

[71] Z. Zhu and R. Soricut. H-transformer-1d: Fast one-dimensional hierarchical attention for sequences, 2021. URL `https://arxiv.org/abs/2107.11906`.

## A  Generalizing Log-Linear Attention to More Expressive Linear RNNs

The main paper adopts the following unified view of efficient attention (Eq. 1):

$$\mathbf{P} = \mathbf{A} \odot \mathbf{M}, \quad \mathbf{O} = \mathbf{P}\mathbf{V},$$

This formulation reveals that the key difference between linear and log-linear attention lies in the structure of the mask matrix $\mathbf{M} \in \mathbb{R}^{T \times T}$. Variations among linear attention models—such as Mamba-2 and Gated DeltaNet—stem from different parameterizations of $\mathbf{A}$. While this perspective offers a unifying and intuitive framework that captures a wide range of attention mechanisms, it comes with an important limitation: the state-transition terms are restricted to be scalars (in the case of Mamba-2) or identity-plus-rank-one matrices (in the case of Gated DeltaNet).

In this section, we introduce a more general framework that relaxes this scalar constraint by allowing state-transition terms (including the thus $\lambda_t^{(\ell)}$ terms) to be matrix-valued. This extension enables richer and more expressive attention mechanisms while preserving computational efficiency.

**Linear Attention as an SSS Tensor.** Consider the standard linear attention mechanism with data-dependent gating and an SSS (sequentially semiseparable) mask $\mathbf{M}^{\mathcal{S}}$:

$$\mathbf{P} = \mathbf{Q}\mathbf{K}^\top \odot \mathbf{M}^{\mathcal{S}}, \quad \mathbf{O} = \mathbf{P}\mathbf{V}.$$

In the main paper, we extend the SSS mask $\mathbf{M}^{\mathcal{S}}$ to a hierarchical form $\mathbf{M}^{\mathcal{H}}$. Notice that in Mamba-2, the resulting matrix $\mathbf{P}$ also inherits the same structural property, with its SSS-rank governed by the hidden dimension $d$:

$$\mathbf{P}_{t,s} = \mathbf{Q}_t \left( \mathbf{C}_t \cdots \mathbf{C}_{s+1} \right) \mathbf{K}_s^\top, \quad \text{where} \quad \mathbf{C}_t = \alpha_t \mathbf{I}.$$

We now define a 4D tensor $\mathbf{M}^{\mathcal{S}} \in \mathbb{R}^{(T \times T) \times (d \times d)}$ such that:

$$\mathbf{P}_{t,s} = \mathbf{Q}_t \mathbf{M}_{t,s} \mathbf{K}_s^\top, \quad \text{where} \quad \mathbf{M}_{t,s} = \mathbf{C}_t \cdots \mathbf{C}_{s+1}.$$

Each entry $\mathbf{M}_{t,s} \in \mathbb{R}^{d \times d}$ is a matrix, making $\mathbf{M}^{\mathcal{S}}$ a 4D tensor. We refer to this as an SSS tensor due to its sequentially semiseparable-like structure along the temporal dimension, though this term is not yet formalized in the literature.
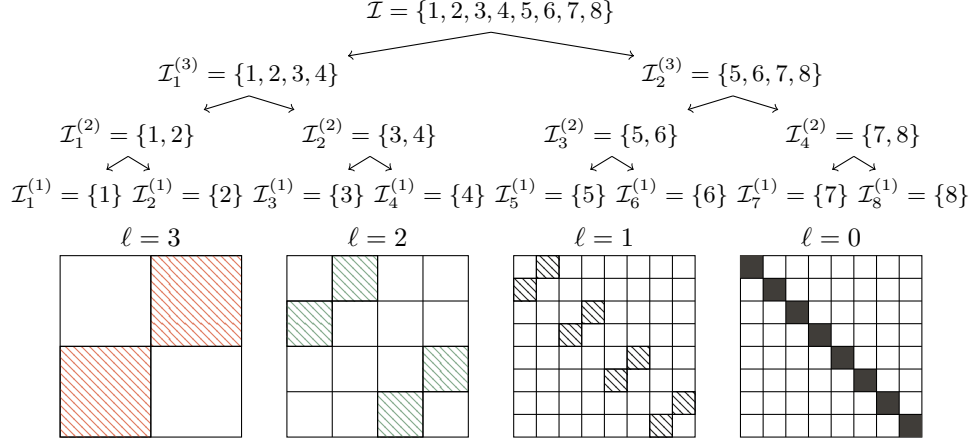
$$\mathcal{I} = \{1, 2, 3, 4, 5, 6, 7, 8\}$$

$$\mathcal{I}_1^{(3)} = \{1, 2, 3, 4\} \qquad\qquad \mathcal{I}_2^{(3)} = \{5, 6, 7, 8\}$$

$$\mathcal{I}_1^{(2)} = \{1, 2\} \qquad \mathcal{I}_2^{(2)} = \{3, 4\} \qquad \mathcal{I}_3^{(2)} = \{5, 6\} \qquad \mathcal{I}_4^{(2)} = \{7, 8\}$$

$$\mathcal{I}_1^{(1)} = \{1\} \;\; \mathcal{I}_2^{(1)} = \{2\} \;\; \mathcal{I}_3^{(1)} = \{3\} \;\; \mathcal{I}_4^{(1)} = \{4\} \;\; \mathcal{I}_5^{(1)} = \{5\} \;\; \mathcal{I}_6^{(1)} = \{6\} \;\; \mathcal{I}_7^{(1)} = \{7\} \;\; \mathcal{I}_8^{(1)} = \{8\}$$

$$\ell = 3 \qquad\qquad \ell = 2 \qquad\qquad \ell = 1 \qquad\qquad \ell = 0$$



**Figure 8:** Visualization adapted from [37, 31]: This example illustrates a cluster tree of depth 3 along with the corresponding block partitions at each level. Blocks marked with stripes are stored as low-rank matrices in the HODLR format, while those filled with solid color represent dense matrices.

This tensor-centric view naturally accommodates matrix-valued state transitions $\mathbf{C}_t \in \mathbb{R}^{d \times d}$ with arbitrary structure, offering a richer representation than scalar- or identity-plus-rank-one-based approaches. In particular, models such as Mamba-2 and Gated DeltaNet can be interpreted as operating on 4D tensors with different hidden-dimension structures, while still preserving temporal semiseparability.[13]

$$\text{Mamba-2:} \quad \mathbf{M}_{t,s}^{\mathcal{S}} = \prod_{t'=t}^{s+1} \alpha_{t'} \mathbf{I}, \qquad\qquad \text{Gated DeltaNet:} \quad \mathbf{M}_{t,s}^{\mathcal{S}} = \prod_{t'=t}^{s+1} \alpha_{t'} \left( \mathbf{I} - \beta_{t'} \boldsymbol{k}_{t'} \boldsymbol{k}_{t'}^{\top} \right)$$

**Log-Linear Attention as an $\mathcal{H}$ Tensor.** We can apply our *log-linear* attention to these more flexible (linear) RNNs by incorporating matrix-valued, level- and data-dependent terms $\boldsymbol{\Lambda}_t^{(\ell)} \in \mathbb{R}^{d \times d}$:

$$\text{Mamba-2:} \quad \mathbf{M}_{t,s}^{\mathcal{H}} = \boldsymbol{\Lambda}_t^{(\ell)} \prod_{t'=t}^{s+1} \alpha_{t'} \mathbf{I}, \qquad\qquad \text{Gated DeltaNet:} \quad \mathbf{M}_{t,s}^{\mathcal{H}} = \boldsymbol{\Lambda}_t^{(\ell)} \prod_{t'=t}^{s+1} \alpha_{t'} \left( \mathbf{I} - \beta_{t'} \boldsymbol{k}_{t'} \boldsymbol{k}_{t'}^{\top} \right)$$

This formulation highlights a key insight: both Mamba-2 and Gated DeltaNet share a common semiseparable structure in the temporal dimension, but differ in how they structure the hidden dimension. Mamba-2 relies on scaled identities, while Gated DeltaNet applies identity-minus-rank-one modifications. Table 5 summarizes these distinctions.

| Model | Temporal Structure | Hidden Size Structure |
|---|---|---|
| Mamba-2 | Semiseparable | Scaled Identity |
| Gated DeltaNet | Semiseparable | Identity plus Low-Rank |
| *Log-Linear* Mamba-2 | Hierarchical | Scaled Identity |
| *Log-Linear* Gated DeltaNet | Hierarchical | Identity plus Low-Rank |

**Table 5:** Structural comparison of different attention variants.

# B  Log-Linear Attention as $\mathcal{H}$ Matrices

We begin by introducing two classes of Hierarchical matrices ($\mathcal{H}$ matrices) following Massei et al. [37]: HODLR (Hierarchically Off-Diagonal Low-Rank) matrices and HSS (Hierarchically Semi-Separable) matrices. We then show how Log-Linear Attention corresponds to a specific subclass of $\mathcal{H}$ matrices that occupies an intermediate position between these two. Finally, we discuss a further variant of $\mathcal{H}$ matrices that, in principle, allows for more refined partitioning—potentially enhancing approximation quality at the cost of increased (though constant-factor) computational complexity.

---

[13]Strictly speaking, Gated DeltaNet also need to include a term $\beta_t$ from $\beta_t \boldsymbol{v}_t \boldsymbol{k}_t^{\top}$. For clarity, we omit it here, as it can be absorbed into other terms.

## B.1 HODLR Matrices

HODLR (Hierarchically Off-Diagonal Low-Rank) matrices are structured matrices built via recursive partitioning, where off-diagonal blocks are low-rank at every level. This structure is formalized using a cluster tree [37]. Let $T$ be the matrix dimension, and let $\mathcal{T}$ be a perfectly balanced binary tree of depth $L$ whose nodes are subsets of $\{1, \ldots, T\}$. We say $\mathcal{T}$ is a cluster tree if: (1) the root is $\mathcal{I} = \{1, \ldots, T\}$; (2) each level partitions indices into contiguous blocks; (3) every node $\mathcal{I}^{(\ell)}i$ at level $\ell$ has two children $\mathcal{I}_{2i-1}^{(\ell-1)}$ and $\mathcal{I}_{2i}^{(\ell-1)}$ that form a disjoint partition of the parent. See Fig. 8 for a visual example of such a hierarchical partitioning.

Now, let $\mathbf{M} \in \mathbb{R}^{T \times T}$ be a square matrix and $\mathcal{T}$ a cluster tree as described above. We say that $\mathbf{M}$ is a $(\mathcal{T}, k)$-HODLR matrix if,

$$\mathrm{rank}\left(\mathbf{M}[\mathcal{I}_i^{(\ell)}, \mathcal{I}_j^{(\ell)}]\right) \leq k, \quad \forall\ \mathcal{I}_i^{(\ell)}, \mathcal{I}_j^{(\ell)} \in \mathrm{siblings}\,(\mathcal{T})$$

This hierarchical low-rank structure enables efficient $\mathcal{O}(T \log T)$ storage and matrix-vector multiplication, making HODLR matrices a core component in fast algorithms for dense matrix computations. HODLR belongs to the broader class of rank-structured matrices known as Hierarchical matrices ($\mathcal{H}$ matrices).

## B.2 HSS Matrices

The $\mathcal{O}(T \log T)$ memory complexity of HODLR matrices arises from their recursive structure: they consist of $\mathcal{O}(\log T)$ levels, each storing low-rank factorizations that require $\mathcal{O}(T)$ space. In cases where these low-rank factors exhibit linear dependencies across levels, it is possible to exploit these relationships through nested hierarchical low-rank representations, potentially reducing the memory complexity to $\mathcal{O}(T)$ by eliminating the logarithmic factor [37].

Let $\mathcal{I}_i^{(\ell)}$ and $\mathcal{I}_j^{(\ell)}$ denote a pair of sibling clusters at level $\ell$ in the cluster tree $\mathcal{T}$. Define $n^{(\ell)} = 2^{\ell-1}$ as the block size at level $\ell$. The off-diagonal block corresponding to these clusters can be parameterized as:

$$\mathbf{M}[\mathcal{I}_i^{(\ell)}, \mathcal{I}_j^{(\ell)}] = \mathbf{U}_i^{(\ell)} \mathbf{\Sigma}_{i,j}^{(\ell)} \left(\mathbf{V}_j^{(\ell)}\right)^\top, \quad \text{where} \quad \mathbf{U}_i^{(\ell)}, \mathbf{V}_j^{(\ell)} \in \mathbb{R}^{n^{(\ell)} \times k}, \ \mathbf{\Sigma}_{i,j}^{(\ell)} \in \mathbb{R}^{k \times k}$$

We call $\mathbf{M}$ matrix a Hierarchically Semiseparable matrices (HSS) if low-rank factors at different levels are linearly related through some "translation operators" $\mathbf{T}_{\mathbf{U}}^{(\ell)}, \mathbf{T}_{\mathbf{V}}^{(\ell)} \in \mathbb{R}^{2k \times k}$ such that,

$$\mathbf{U}_i^{(\ell)} = \begin{bmatrix} \mathbf{U}_{i_1}^{(\ell-1)} & 0 \\ 0 & \mathbf{U}_{i_2}^{(\ell-1)} \end{bmatrix} \mathbf{T}_{\mathbf{U},i}^{(\ell)}, \quad \mathbf{V}_j^{(\ell)} = \begin{bmatrix} \mathbf{V}_{j_1}^{(\ell-1)} & 0 \\ 0 & \mathbf{V}_{j_2}^{(\ell-1)} \end{bmatrix} \mathbf{T}_{\mathbf{V},j}^{(\ell)}$$

More broadly, HSS matrices belong to a subclass of $\mathcal{H}$ matrices known as $\mathcal{H}^2$ matrices.

## B.3 Quasi-Hierarchical Matrix.

As discussed above, when the low-rank basis matrices $\mathbf{U}^{(\ell)}$ and $\mathbf{V}^{(\ell)}$ exhibit linear relationships across levels $\ell$, the matrix $\mathbf{M}$ reduces to a semiseparable form. In this case, both storage and matrix-vector multiplication complexities can be reduced to $\mathcal{O}(T)$. Otherwise, $\mathbf{M}$ retains the general hierarchical structure with $\mathcal{O}(T \log T)$ complexity.

We define a *Quasi-Hierarchical Matrix* as one in which only one of the basis sequences, either $\mathbf{U}^{(\ell)}$ or $\mathbf{V}^{(\ell)}$, satisfies such a linear nesting property across levels, while the other does not. The matrix $\mathbf{M}^{\mathcal{H}}$ used in the Log-Linear model (Eq. 4) is an instance of this structure.

Both Hierarchical and Quasi-Hierarchical matrices incur $\mathcal{O}(T \log T)$ complexity for storage and computation during training. However, the use of Quasi-Hierarchical matrices plays a crucial role in enabling $\mathcal{O}(\log T)$ complexity during inference. We are not aware of a recurrent algorithm for general Hierarchical matrices that achieves logarithmic inference complexity.[14]

---

[14]In fact, our initial attempts involved using fully Hierarchical matrices, but we were unable to derive a recurrent formulation with $\mathcal{O}(\log T)$ complexity. This motivated the design of Quasi-Hierarchical matrices specifically to support efficient recurrence.
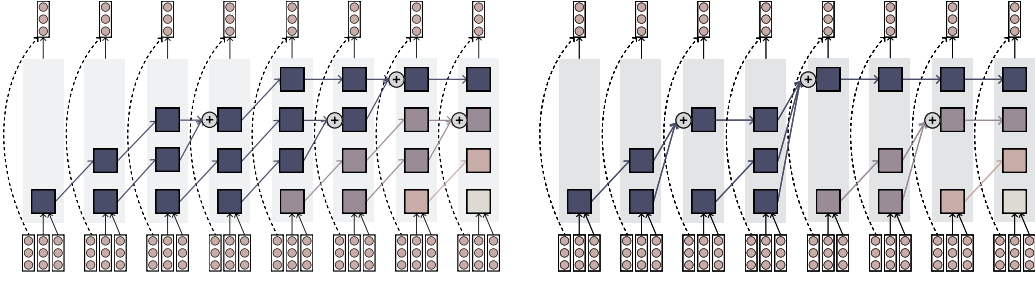
Figure 9: **Left**: $\mathcal{H}$ matrices with strong admissibility. **Right**: $\mathcal{H}$ matrices with weak admissibility.

**Reparameterization.** More precisely, Eq. 4 represents a Quasi-Hierarchical matrix that has been specifically re-parameterized as a composition of the scalar weights $\lambda^{(\ell)}$ and a sequentially semisep-arable (SSS) matrix $\mathbf{M}^{\mathcal{S}}$. This reparameterization serves two purposes: first, to highlight the connection between our use of $\mathcal{H}$ matrices and the SSS format adopted in prior work; and second, to enable the block decomposition into a hierarchy of SSS matrices, as shown in Eq. 3.2.

We present this re-parameterization below, along with its 4D tensor variant discussed in §A, where we additionally assume that the matrices $\mathbf{U}_i$ and $\mathbf{V}_j$ are invertible.

**Matrix:**

$$\mathbf{M}_{i,j}^{\mathcal{H}} := \tau_i^{(\ell)} u_i v_j \Leftrightarrow \lambda_i^{(\ell)} \prod_{t=j+1}^{i} \alpha_t$$

$$\Rightarrow \quad \tau_i^{(\ell)} := \lambda_i^{(\ell)}, \ u_i := \prod_{t=0}^{i} \alpha_t, \ v_j := \prod_{t=0}^{j} \frac{1}{\alpha_t}$$

$$\Leftarrow \quad \lambda_i^{(\ell)} := \tau_i^{(\ell)} u_i v_i, \quad a_t := \frac{r_{t-1}}{r_t}$$

**Tensor:**

$$\mathbf{M}_{i,j}^{\mathcal{H}} := \mathbf{T}_i^{(\ell)} \mathbf{U}_i \mathbf{V}_j^{\top} \Leftrightarrow \mathbf{\Lambda}_i^{(\ell)} \prod_{t=i}^{j+1} \mathbf{C}_t$$

$$\mathbf{T}_i^{(\ell)} := \mathbf{\Lambda}_i^{(\ell)}, \ \mathbf{U}_i := \prod_{t=i}^{0} \mathbf{C}_t, \ \mathbf{V}_j^{\top} := \prod_{t=0}^{j} \mathbf{C}_t^{-1}$$

$$\mathbf{\Lambda}_i^{(\ell)} := \mathbf{T}_i^{(\ell)} \mathbf{U}_i \mathbf{V}_i^{\top}, \quad \mathbf{C}_t := \mathbf{R}_t^{-1} \mathbf{R}_{t-1}$$

## B.4  $\mathcal{H}$ Matrices with Strong and Weak Admissibility

In the recurrent formulation of Log-Linear Attention, although there are $\mathcal{O}(\log T)$ states correspond-ing to different hierarchical levels, roughly half of them are zero in practice. This sparsity arises from the specific structure of HODLR matrices, which belong to a broader class of $\mathcal{H}$ matrices known as *weakly admissible* [21].

Figures 10 and 9 illustrate an alternative structure based on strong (or standard) admissibility. Unlike the weakly admissible variant, strongly admissible $\mathcal{H}$ matrices allow for finer-grained partitioning of the matrix, and their corresponding recurrent forms utilize all hierarchical levels.

While strong admissibility can yield more accurate approximations, it comes with a significant com-putational cost [21]. In our early experiments, using strong admissibility in a `Triton` implementation resulted in up to a 4x slowdown, with only marginal improvements in accuracy. As a result, we adopt the weakly admissible structure throughout this work and refer to it simply as the $\mathcal{H}$-matrix.
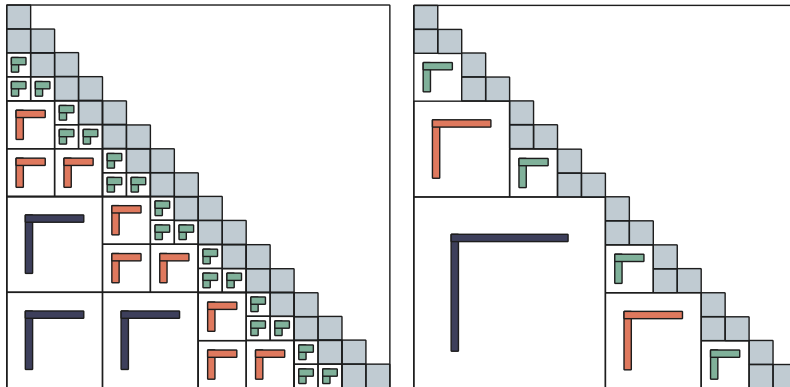


Figure 10: **Left**: $\mathcal{H}$ matrices with strong admissibility. **Right**: $\mathcal{H}$ matrices with weak admissibility.

18

## C   Implementations

```python
import torch
import numpy as np
import torch.nn.functional as F


def segsum(x):
    T = x.size(-1)
    x_cumsum = torch.cumsum(x, dim=-1)
    x_segsum = x_cumsum[..., :, None] - x_cumsum[..., None, :]
    mask = torch.tril(torch.ones(T, T, device=x.device, dtype=bool))
    x_segsum = x_segsum.masked_fill(~mask, -torch.inf)
    return x_segsum


def level_mask(level, T):
    if level == 0:
        return torch.eye(T, dtype=torch.bool)

    i, j = torch.meshgrid(torch.arange(T), torch.arange(T), indexing="ij")
    half = 1 << (level - 1)
    clipped = i - (i %
    valid = (i %
    return valid


def construct_H_matrix(a, L):
    T = a.size(-1)
    A = torch.exp(segsum(a))
    return sum([A * L[..., level, :].unsqueeze(-1) * level_mask(level, T) for level in range(int(np.
     log2(T)) + 1)])


def hattention(X, A, B, C, L, block_len=8):
    """
    Arguments:
    X: (batch, length, n_heads, d_head)
    A: (batch, length, n_heads)
    B: (batch, length, n_heads, d_state)
    C: (batch, length, n_heads, d_state)
    L: (batch, length, n_heads, num_levels) where num_levels = log2(length) + 1
    Return:
    Y: (batch, length, n_heads, d_head)
    """
    T = X.shape[1]
    assert X.dtype == A.dtype == B.dtype == C.dtype
    assert X.shape[1] %
    input_shape = X.shape
    # Rearrange into blocks/chunks
    b, cl = X.shape[0], X.shape[1]
    c = cl // block_len
    X, A, B, C, L = [x.reshape(b, c, block_len, *x.shape[2:]) for x in (X, A, B, C, L)]
    A = A.permute(0, 3, 1, 2)   # (batch, n_heads, c, block_len)
    A_cumsum = torch.cumsum(A, dim=-1)   # (batch, n_heads, c, block_len)

    num_intra_chunk_levels = int(np.log2(block_len)) + 1
    num_inter_chunk_levels = int(np.log2(T)) + 1 - num_intra_chunk_levels
    # Partition the lambda into intra-chunk and inter-chunk lambda
    L_intra, L_inter = L[..., :num_intra_chunk_levels], L[..., num_intra_chunk_levels:]
    L_intra = L_intra.permute(0, 3, 1, 4, 2)   # (batch, n_heads, num_chunks, num_levels, block_len)

    # Intra-chunk Computation
    H = construct_H_matrix(A, L_intra)   # Materialize the H matrix as a dense matrix
    Y_diag = torch.einsum("bclhn,bcshn,bhcls,bcshp->bclhp", C, B, H, X)

    # Inter-chunk Computation
    decay_states = torch.exp((A_cumsum[..., -1:] - A_cumsum))
    states = torch.einsum("bclhn,bhcl,bclhp->bchpn", B, decay_states, X)
    decay_chunk = F.pad(torch.exp(segsum(A_cumsum[..., -1])), (0, 0, 1, 0))[..., :-1, :]
    state_decay_out = torch.exp(A_cumsum)

    def compute_Y_off_level(states, level):
        mask = level_mask(level + 1, c).unsqueeze(0).unsqueeze(0)
        decay_chunk_level = decay_chunk * mask
        states = torch.einsum("bhzc,bchpn->bzhpn", decay_chunk_level, states)
        Y_off = torch.einsum(
            "bclhn,bchpn,bhcl,bclh->bclhp",
            C,
            states,
            state_decay_out,
            L_inter[..., level],
        )
        return Y_off

    Y_off = torch.zeros_like(Y_diag)
    for i in range(num_inter_chunk_levels):
        Y_off += compute_Y_off_level(states, i)

    Y = (Y_off + Y_diag).reshape(input_shape)
    return Y
```

---

**Algorithm 1** Chunkwise Log-Linear Attention Algorithm

---

1: **for** $t \in [T/C]$ **do**
2: $\quad \mathbf{Y}_{[t]} = \left( \mathbf{Q}_{[t]} \mathbf{K}_{[t]}^{\top} \odot \mathbf{M}_{[t]}^{\mathcal{H}} \right) \mathbf{V}_{[t]}$
3: **end for**
4:
5: **for** $\ell \in [\log_2 (T/C)]$ **do**
6: $\quad$ **for** $t \in [T/C]$ **do**
7: $\quad\quad \mathbf{Y}_{[t]} = \mathbf{Y}_{[t]} + \text{mask}_{\mathbf{Q}}^{(\ell)} \left( \mathbf{\Lambda}_{[t]}^{(\ell)} \odot \mathbf{Q}_{[t]} \mathbf{S}_{[t]} \right)$
8: $\quad\quad \mathbf{S}_{[t+1]} = \text{mask}_{\mathbf{A}}^{(\ell)} \left( \mathbf{A}_{[t]} \mathbf{S}_{[t]} \right) + \text{mask}_{\mathbf{K}}^{(\ell)} \left( \mathbf{K}_{[t]} \mathbf{V}_{[t]}^{\top} \right)$
9: $\quad$ **end for**
10: **end for**
11: **return** $\mathbf{Y}$

---

A naive implementation computes each level independently using a Mamba-2-style primitive, then sums the outputs—leading to redundant memory access and kernel launches. To improve efficiency, we fuse computation across four levels into a single Triton kernel, which we found optimal given SRAM constraints on an H100.

For backpropagation, we unify gradient computation across all levels for $\nabla \mathbf{K}$ and $\nabla \mathbf{V}$ by analytically factoring their dependencies. This reduces kernel count and improves memory efficiency, achieving over 3× speedup compared to the naive multi-level version.

## D  Additional Experiment Details

For the implementation benchmarks, all experiments were conducted on an H100 GPU with a batch size of 2, using 48 attention heads, a head dimension of 64, and a chunk size of 64. In Mamba-2-style models, the attention heads are applied to $\mathbf{V}$ (MVA pattern), whereas in FlashAttention-2, we adopt GQA-style attention by applying heads to $\mathbf{Q}$. The dimensions of the $\mathbf{Q}$ and $\mathbf{K}$ states are set to 128, aligning with common training configurations.

For the language modeling experiments, each run was performed on $8\times$A100 or $8\times$H100 GPUs over the course of several days. We do not tie word embeddings, use a vocabulary size of 32,000, and set the initializer range to 0.006. Training is performed with a global batch size of approximately 524K tokens for 95K steps (roughly 50B tokens). We use the `flash-linear-attention` and `flame` libraries [64, 70], following most of their default configurations.

**Detailed Experimental Results.** Table 6 provide detailed results on the Needle-In-A-Haystack task.

| | S-NIAH-1 (pass-key retrieval) | | | S-NIAH-2 (number in haystack) | | | S-NIAH-3 (uuid in haystack) | | |
|---|---|---|---|---|---|---|---|---|---|
| Model | 4K | 8K | 16K | 4K | 8K | 16K | 4K | 8K | 16K |
| Transformer | 72.6 | 76.0 | 16.6 | 100.0 | 99.8 | 90.0 | 77.4 | 67.0 | 44.6 |
| w/ *24 Layers* | 92.4 | 78.4 | 89.8 | 100.0 | 100.0 | 99.6 | 84.0 | 63.6 | 36.4 |
| Mamba-2 | 90.4 | 56.8 | 21.6 | 72.4 | 28.0 | 18.6 | 4.0 | 3.6 | 0.8 |
| w/ *Log-Linear* | 100.0 | 99.8 | 72.4 | 89.8 | 68.2 | 12.8 | 33.6 | 22.6 | 2.0 |
| Gated DeltaNet | 100.0 | 100.0 | 100.0 | 95.8 | 46.8 | 5.0 | 66.2 | 14.6 | 6.0 |
| w/ *Log-Linear* | 100.0 | 100.0 | 100.0 | 95.6 | 59.6 | 9.2 | 48.8 | 13.0 | 8.8 |

| | MK-NIAH-1 (multi-key line retrieval) | | | MQ-NIAH (multi-query) | | | MV-NIAH (multi-value) | | |
|---|---|---|---|---|---|---|---|---|---|
| Model | 4K | 8K | 16K | 4K | 8K | 16K | 4K | 8K | 16K |
| Transformer (L) | 79.4 | 83.0 | 61.4 | 58.9 | 48.0 | 29.8 | 37.5 | 34.1 | 21.5 |
| w/ *24 Layers* | 62.6 | 83.2 | 75.2 | 54.6 | 46.0 | 34.5 | 48.4 | 45.4 | 32.3 |
| Mamba-2 | 27.2 | 18.6 | 13.6 | 28.7 | 19.4 | 1.3 | 27.9 | 14.8 | 4.4 |
| w/ *Log-Linear* | 43.2 | 39.8 | 21.2 | 26.6 | 22.4 | 6.6 | 28.1 | 22.8 | 8.9 |
| Gated DeltaNet | 23.0 | 21.2 | 5.2 | 21.6 | 16.9 | 7.2 | 16.2 | 14.5 | 7.0 |
| w/ *Log-Linear* | 49.4 | 27.8 | 10.2 | 34.9 | 22.0 | 9.8 | 31.4 | 25.0 | 13.3 |

**Table 6:** NIAH experiments, including three single-needle tasks—S-NIAH-1 (passkey retrieval), S-NIAH-2 (numerical needle), and S-NIAH-3 (UUID-based needle)—and three multi-needle variants: MK-NIAH-1 (multi-key line retrieval), MQ-NIAH (multi-query), and MV-NIAH (multi-value).