

# MAMBA-3: IMPROVED SEQUENCE MODELING USING STATE SPACE PRINCIPLES

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

The recent scaling of test-time compute for LLMs has restricted the practical deployment of models to those with strong capabilities that can generate high-quality outputs in an inference-efficient manner. While current Transformer-based models are the standard, their quadratic compute and linear memory bottlenecks have spurred the development of sub-quadratic models with linear-scaling compute with constant memory requirements. However, many recent linear-style models lack certain capabilities or lag behind in quality, and even their linear-time inference is not hardware-efficient. Guided by an inference-first perspective, we introduce three core methodological improvements inspired by the state-space model viewpoint of linear models. We combine a: 1) more expressive recurrence, 2) complex state update rule that enables richer state tracking, and 3) multi-input, multi-output formulation together, resulting in a stronger model that better exploits hardware parallelism during decoding. Together with architectural refinements, our **Mamba-3** model achieves significant gains across retrieval, state-tracking, and downstream language modeling tasks. Our new architecture sets the Pareto-frontier for performance under a fixed inference budget and outperforms strong baselines in a head-to-head comparison.

## 1 INTRODUCTION

Test-time compute has emerged as a key driver of progress in AI, with techniques like chain-of-thought reasoning and iterative refinement demonstrating that inference-time scaling can unlock new capabilities (Wu et al., 2025; Snell et al., 2024). This paradigm shift makes inference efficiency (Kwon et al., 2023; Li et al., 2024) paramount, as the practical impact of AI systems now depends critically on their ability to perform large-scale inference during deployment. Model architecture design plays a fundamental role in determining inference efficiency, as architectural choices directly dictate the computational and memory requirements during generation. While Transformer-based models (Vaswani et al., 2017) are the current industry standard, they are fundamentally bottlenecked by linearly increasing memory demands through the KV cache and quadratically increasing compute requirements through the self-attention mechanism. These drawbacks have motivated recent lines of work on sub-quadratic models, e.g., state-space models (SSMs), which, despite utilizing only constant memory and linear compute, have comparable or better performance than their Transformer counterparts. Models that benefit the most from this new scaling paradigm perform well on the following three axes: (i) quality, (ii) capability, and (iii) inference efficiency.

Recent model architectures have tried to strike a balance between the three, but many fall short on at least one of these three axes. In particular, Mamba-2 and Gated DeltaNet, which have gained significant traction and adoption due to their inference efficiency, made architectural design choices that enable their linear compute requirements but sacrifice quality and capabilities (Dao & Gu, 2024; Yang et al., 2025a). For example, Mamba-2 was developed to improve training speed and simplicity over Mamba-1 (Gu & Dao, 2024), opting out of more expressive parameterizations of the underlying SSM and hindering the quality of the model (Dao & Gu, 2024). Linear attention-style models (Katharopoulos et al., 2020) have also been shown to lack certain capabilities, with poor state-tracking abilities, e.g., determining parity of bit sequences, being one of the most notable (Grazzi et al., 2025; Sarrof et al., 2024). In addition, despite these sub-quadratic models retaining linear inference, their inference itself is not hardware efficient. Because these algorithms were developed from a training perspective, their decoding phase has low arithmetic intensity (the ratio of FLOPs to memory traffic), resulting in large portions of hardware remaining idle.

To develop more performant models from an inference-first paradigm, in this paper we introduce three core methodological changes, influenced by a SSM-centric viewpoint of sub-quadratic models, on top of Mamba-2. While many recent models fall into the linear attention framework (Dao & Gu, 2024; Yang et al.,

2025a; Sun et al., 2023), we find that the classical SSM toolbox (Kalman, 1960; Gopal, 1993) leads to natural interpretations and improvements on modeling.

**Trapezoidal Discretization.** We discretize the underlying continuous-time dynamical system with a trapezoidal methodology. The final recurrence is a more expressive superset of Mamba-2’s recurrence and can be viewed as a convolution. We combine this new discretization with applied biases on the  $B, C$ , inspired by Yu & Erichson (2025), and find that their synergy is able to empirically replace the short causal convolution in language modeling.

**Complexified State-Space Model.** By viewing the underlying SSM of Mamba-3 as complex-valued, we enable a more expressive state update compared to Mamba-2. This change in update rule, designed to be lightweight for training and inference, overcomes the lack of state-tracking ability common for many current linear models. We highlight that our complex-valued update rule is equivalent to a data-dependent rotary embedding and thus can be calculated efficiently (Su et al., 2023).

**Multi-Input, Multi-Output SSM.** To improve FLOP-efficiency during decoding, we shift from outer-product-based state update to matrix-multiplication-based state update. In view of the signal processing foundations of SSMs, such a transition exactly coincides with the generalization from a single-input single-output (SISO) sequence dynamic to a multiple-input multiple-output (MIMO) one. Here, we found that MIMO is particularly suitable for inference, as the extra expressivity allows for more compute during state update, without increasing the state size and hence compromising speed.

These three SSM-centric methodological changes are core to our **Mamba-3** mixer primitive. We also make adjustments to the overall architecture to ensure more similarity to the baseline Transformer architecture. Mamba-3 swaps the pre-output projection norm with the more common QK-normalization (Team et al., 2025; OLMo et al., 2025) and makes the short convolution, a common component found in many other sub-quadratic models (Gu & Dao, 2024; Yang et al., 2025a; von Oswald et al., 2025), optional.

We empirically validate our new model on a suite of synthetic and language-modeling tasks.

- **Better Quality.** Mamba-3 matches or outperforms Mamba-2 and other open-source architectures on standard downstream language modeling evaluations. For example, Mamba-3-1.5B’s average accuracy on all downstream tasks is better than that of its Transformer, Mamba-2, and Gated DeltaNet counterparts.
- **Better Capability.** Mamba-3’s complexification of the SSM state enables the model to solve synthetic state-tracking tasks that Mamba-2 cannot. We empirically demonstrate that the efficient RoPE-like calculation is able to near perfectly solve arithmetic tasks, while Mamba-3 without RoPE and Mamba-2 perform not better than random guessing.
- **Better Inference Efficiency.** Mamba-3’s MIMO variant retains the same state size while enabling better hardware utilization compared to standard Mamba-3 and other models. Its improved performance without increasing memory requirements pushes the pareto-frontier of inference efficiency.

## 2 PRELIMINARIES

### 2.1 NOTATION

Scalars are denoted by plain-text letters (e.g.,  $x, y$ ). Tensors, including vectors and matrices, are denoted by bold letters (e.g.,  $\mathbf{h}, \mathbf{C}$ ). The shape of the tensor can be inferred from the context. We denote the input sequence length as  $T$ , the model dimension as  $D$ , and the SSM state size as  $N$ . For time indices, we use subscripts (e.g.,  $x_t$  for the input at time  $t$ ). The Hadamard product between two tensors is denoted by  $\odot$ . For a vector of size  $\mathbf{v} \in \mathbb{R}^d$ , we denote  $\text{Diag}(\mathbf{v}) \in \mathbb{R}^{d \times d}$  as the diagonal matrix with the vector  $\mathbf{v}$  as the diagonal, and for products of scalars across time steps, we use the notation  $\alpha_{t \dots s} = \alpha_{t:s}^\times = \prod_{i=s}^t \alpha_i$ .

### 2.2 SSM PRELIMINARIES

State Space Models (SSMs) describe continuous-time linear dynamics via

$$\dot{\mathbf{h}}(t) = \mathbf{A}(t)\mathbf{h}(t) + \mathbf{B}(t)x(t), \quad y(t) = \mathbf{C}(t)^\top \mathbf{h}(t),$$

where  $\mathbf{h}(t) \in \mathbb{R}^N$  is the hidden state,  $x(t) \in \mathbb{R}$  the input, and  $\mathbf{A}(t) \in \mathbb{R}^{N \times N}$ ,  $\mathbf{B}(t), \mathbf{C}(t) \in \mathbb{R}^N$ . For discrete sequences with step size  $\Delta_t$ , Euler’s discretization gives the recurrence

$$\mathbf{h}_t = e^{\Delta_t \mathbf{A}_t} \mathbf{h}_{t-1} + \Delta_t \mathbf{B}_t x_t, \quad y_t = \mathbf{C}_t^\top \mathbf{h}_t.$$

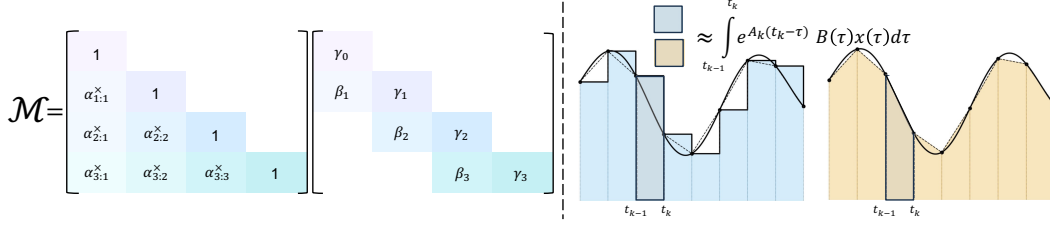


Figure 1: **Left:** The structured mask induced by the generalized trapezoid rule is a product of the decay and convolutional mask. **Right:** Euler (hold endpoint) vs trapezoidal rule (average endpoints).

**Mamba-2’s parameterization.** Mamba-2 (Dao & Gu, 2024) makes the SSM *data-dependent* and hardware-efficient by (i) projecting  $A = \mathbf{A} \in \mathbb{R}_{<0}$ , and  $\mathbf{B}, \mathbf{C} \in \mathbb{R}^N$  from the current token and (ii) choosing transition matrix  $A = \mathbf{A}$  as a data-dependent scalar. Writing  $\alpha_t := e^{\Delta_t A_t} \in (0, 1)$  and  $\gamma_t := \Delta_t$ , the update becomes

$$\mathbf{h}_t = \alpha_t \mathbf{h}_{t-1} + \gamma_t \mathbf{B}_t x_t, \quad y_t = \mathbf{C}_t^\top \mathbf{h}_t.$$

The scalar  $A_t < 0$  is an input-dependent *forget-gate (decay)*  $\alpha_t$ , and the parameter *selectivity*  $\Delta_t$  jointly controls the forget-gate ( $\alpha_t = \exp(\Delta_t A_t)$ ) and the input-gate ( $\gamma_t = \Delta_t$ ): larger  $\Delta_t$  forgets faster and up-weights the current token more strongly, while smaller  $\Delta_t$  retains the hidden state with minimal contributions from the current token.

### 2.3 STRUCTURED MASKED REPRESENTATION AND STATE SPACE DUALITY

Dao & Gu (2024) show that a large class of SSMs admit a *matrix* form that vectorizes the time-step recurrence. For instance, Mamba-2’s recurrence can be vectorized as a masked matrix multiplication,

$$\mathbf{Y} = (\mathbf{L} \odot \mathbf{C} \mathbf{B}^\top) \mathbf{X} = \left( \begin{bmatrix} 1 & & & \\ \alpha_1 & 1 & & \\ \vdots & & \ddots & \\ \alpha_{T-1} & \dots & \alpha_T & 1 \end{bmatrix} \odot \mathbf{C} \mathbf{B}^\top \right) \mathbf{X}, \quad (1)$$

where  $\mathbf{L} \in \mathbb{R}^{T \times T}$  is the structured mask,  $\mathbf{B}, \mathbf{C} \in \mathbb{R}^{T \times N}$ ,  $\mathbf{X} \in \mathbb{R}^{T \times D}$  is the input to the SSM and  $\mathbf{Y} \in \mathbb{R}^{T \times D}$  is its output. Within this form, Mamba-2 can be viewed as a type of linear attention by setting  $\mathbf{Q} = \mathbf{C}$ ,  $\mathbf{K} = \mathbf{B}$ ,  $\mathbf{V} = \mathbf{X}$  and viewing  $\mathbf{L}$  as a causal, data-dependent mask. When all  $\alpha = 1$ , the expression reduces to (causal) linear attention (Katharopoulos et al., 2020). A more detailed coverage of related linear-time sequence mixers can be found at Appendix A.

## 3 MODEL DESIGN FROM A STATE-SPACE VIEWPOINT

We introduce Mamba-3, with three new innovations rooted in classical state-space theory: trapezoidal discretization for more expressive dynamics, complex-valued state spaces for state-tracking, and multi-input multi-output (MIMO) to improve hardware utilization. These advances address the quality, capability, and efficiency limitations of current sub-quadratic architectures.

### 3.1 TRAPEZOIDAL DISCRETIZATION

Structured SSMs are naturally defined as continuous-time dynamical systems that map input functions,  $x(t) \in \mathbb{R}$ , to output functions,  $y(t) \in \mathbb{R}$ , for time  $t > 0$ . In sequence modeling, however, the data is only observed at discrete time steps, which then requires applying a *discretization step* to the SSM to transform its continuous-time dynamics into a discrete recurrence. The preliminary step in deriving Mamba-3’s discretization is to apply the Variation of Constants formula (Proposition 5), which decomposes the hidden state into an exponentially decay term and a state update term ‘information’ term dependent on the most recent inputs.

The first step in deriving the discretized recurrence is to approximate the “state-update” integral in equation 10. A straightforward choice, used in Mamba-2, is applying *Euler’s rule* (Süli & Mayers, 2003), which approximates the integral by holding the (right) endpoint constant throughout the interval (Fig. 1). This yields Mamba-2’s recurrence,

$$\begin{aligned} \mathbf{h}_t &= e^{\Delta_t A_t} \mathbf{h}_{t-1} + (\tau_t - \tau_{t-1}) e^{(\tau_t - \tau_{t-1}) A_t} \mathbf{B}_t x_t \\ &\approx e^{\Delta_t A_t} \mathbf{h}_{t-1} + \Delta_t \mathbf{B}_t x_t. \end{aligned} \quad (2)$$

However, Euler’s rule provides only a first-order approximation to the “state-update” integral: local truncation error is  $O(\Delta_t^2)$ , which accumulates across steps to yield a global error of  $O(\Delta_t)$  over the sequence. In contrast, we adopt a *generalized trapezoidal rule*, which provides a second-order accurate approximation of the integral, offering improved accuracy over the Euler’s rule. Specifically, it approximates the integral with a *data-dependent, convex combination of both interval endpoints*. This generalization extends the classical

trapezoidal rule (Süli & Mayers, 2003), which simply averages the interval endpoints, by allowing for a *data-dependent convex combination* (Fig. 1).

**Proposition 1** (Generalized Trapezoidal Discretization). *Approximating the state-update integral in equation 10 by the general trapezoidal rule yields the recurrence,*

$$\mathbf{h}_t = e^{\Delta_t A_t} \mathbf{h}_{t-1} + (1 - \lambda_t) \Delta_t e^{\Delta_t A_t} \mathbf{B}_{t-1} x_{t-1} + \lambda_t \Delta_t \mathbf{B}_t x_t, \quad (3)$$

$$:= \alpha_t \mathbf{h}_{t-1} + \beta_t \mathbf{B}_{t-1} x_{t-1} + \gamma_t \mathbf{B}_t x_t, \quad (4)$$

where  $\lambda_t \in [0, 1]$  is a data-dependent scalar,  $\alpha_t := e^{\Delta_t A_t}$ ,  $\beta_t := (1 - \lambda_t) \Delta_t e^{\Delta_t A_t}$ ,  $\gamma_t := \lambda_t \Delta_t$ .

**Remark 1** (Expressivity). Our scheme is a generalization of a) The classical trapezoid rule which is recovered when  $\lambda_t = \frac{1}{2}$ . b) Mamba-2’s Euler’s rule, which is recovered when  $\lambda_t = 1$ .

**Remark 2** (Error Rate). This is a second-order discretization with local truncation error  $O(\Delta_t^3)$  and global error  $O(\Delta_t^2)$  over the sequence under standard stability assumptions.

### 3.1.1 TRAPEZOIDAL DISCRETIZATION IS A CONVOLUTIONAL MASK

We can view the generalized trapezoidal discretization as applying a *data-dependent* convolution of size two on the projected input,  $\mathbf{B}_t x_t$ , to the SSM. We now show that a similar vectorization to Equation (1) holds with the generalized trapezoidal discretization. Unrolling the recurrence starting from  $\mathbf{h}_0 = \gamma_0 \mathbf{B}_0 x_0$  results in  $\mathbf{h}_T = \alpha_T \dots \alpha_2 (\gamma_0 \alpha_1 + \beta_1) \mathbf{B}_0 x_0 + \dots + \gamma_T \mathbf{B}_T x_T$ .

Unrolling these rows shows that the mask induced by the trapezoidal update is no longer a fixed averaging of endpoints (as in the classical trapezoidal rule), but a *data-dependent convex combination* of the two interval endpoints. In the SSD representation, this corresponds to a mask  $\mathbf{L}$ :

$$\begin{bmatrix} \gamma_0 & & & & \\ (\gamma_0 \alpha_1 + \beta_1) & & & & \\ \alpha_2 (\gamma_0 \alpha_1 + \beta_1) & \gamma_2 & & & \\ \vdots & & \ddots & & \\ \alpha_{T \dots 2} (\gamma_0 \alpha_1 + \beta_1) & & \dots & \gamma_T \end{bmatrix} = \begin{bmatrix} 1 & & & & \\ \alpha_1 & 1 & & & \\ \alpha_2 \alpha_1 & & & & \\ \vdots & & \ddots & & \\ \alpha_{T \dots 1} & & \dots & 1 \end{bmatrix} \begin{bmatrix} \gamma_0 & & & & \\ \beta_1 & & & & \\ 0 & \gamma_2 & & & \\ \vdots & & \ddots & & \\ 0 & & \dots & \gamma_T \end{bmatrix}. \quad (5)$$

Here, the first factor is precisely the lower-triangular decay mask from Mamba-2, while the second factor encodes the size two convolution induced by the trapezoidal rule through the coefficients  $(\beta_t, \gamma_t)$ . We provide a rigorous proof for this decomposition in Appendix B.2.

### 3.2 COMPLEX-VALUED SSMs

Modern SSMs are designed with efficiency as the central goal, motivated by the need to scale to larger models and longer sequences. For instance, successive architectures have progressively simplified the state transition matrix: S4 (Gu et al., 2022a) used complex-valued Normal plus Low Rank (NPLR) matrices, Mamba (Gu & Dao, 2024) reduced this to a diagonal of reals, and Mamba-2 (Dao & Gu, 2024) further simplified it to a single scalar. Although these simplifications largely maintain language modeling performance, recent works (Merrill et al., 2025; Sarrof et al., 2024; Grazzi et al., 2025) have shown that they degrade the capabilities of the model on simple state-tracking tasks such as parity and modular arithmetic, which can be solved by a one-layer LSTM.

This limitation, formalized in Theorem-1 of (Grazzi et al., 2024), arises from restricting the eigenvalues of the transition matrix to real numbers, which cannot represent “rotational” hidden state dynamics. For instance, consider the parity function on binary inputs  $\{0, 1\}$ , defined as  $\sum_t x_t \bmod 2$ . This task can be performed using update:  $\mathbf{h}_t = \mathbf{R}(\pi x_t) \mathbf{h}_{t-1}$ , where  $\mathbf{R}(\cdot)$  is a 2-D rotation matrix. Such rotational dynamics cannot be expressed with real eigenvalues.

To recover this capability, we begin with complex SSMs (6), which are capable of representing state-tracking dynamics. We show that, under discretization (Proposition 5), complex SSMs can be formulated as a real SSMs with a *block-diagonal transition matrix composed of  $2 \times 2$  rotation matrices* (Proposition 2). We then show that this is equivalent to applying *data-dependent rotary embeddings* on both the input and output projections  $\mathbf{B}, \mathbf{C}$  respectively. This result establishes a theoretical connection between complex SSMs and data-dependent RoPE embeddings (Proposition 3). Finally, this allows for an efficient implementation of the SSM via “RoPE trick”, enabling efficient complex-valued state transition matrix with minimal overhead over real SSMs.

**Proposition 2** (Complex-to-Real SSM Equivalence). *Consider a complex-valued SSM*

$$\dot{\mathbf{h}}(t) = \text{Diag}(A(t) + i\theta(t)) \mathbf{h}(t) + (\mathbf{B}(t) + i\hat{\mathbf{B}}(t)) x(t), \quad (6)$$

$$y(t) = \text{Re}\left((\mathbf{C}(t) + i\hat{\mathbf{C}}(t))^\top \mathbf{h}(t)\right),$$

where  $\mathbf{h}(t) \in \mathbb{C}^{N/2}$ ,  $\boldsymbol{\theta}(t), \mathbf{B}(t), \hat{\mathbf{B}}(t), \mathbf{C}(t), \hat{\mathbf{C}}(t) \in \mathbb{R}^{N/2}$ , and  $x(t), A(t) \in \mathbb{R}$ . Under Euler discretization, this system is equivalent to a real-valued SSM

$$\begin{aligned} \mathbf{h}_t &= e^{\Delta_t A_t} \mathbf{R}_t \mathbf{h}_{t-1} + \Delta_t \mathbf{B}_t x_t, \\ y_t &= \mathbf{C}_t^\top \mathbf{h}_t, \end{aligned} \quad (7)$$

with state  $\mathbf{h}_t \in \mathbb{R}^N$ , projections

$$\mathbf{B}_t = \begin{bmatrix} \mathbf{B}_t \\ \hat{\mathbf{B}}_t \end{bmatrix} \in \mathbb{R}^N, \quad \mathbf{C}_t = \begin{bmatrix} \mathbf{C}_t \\ -\hat{\mathbf{C}}_t \end{bmatrix} \in \mathbb{R}^N,$$

and a transition matrix

$$\mathbf{R}_t = \text{Block}\left(\{R(\Delta_t \boldsymbol{\theta}_t[i])\}_{i=1}^{N/2}\right) \in \mathbb{R}^{N \times N}, \quad R(\Theta) = \begin{bmatrix} \cos(\Theta) & -\sin(\Theta) \\ \sin(\Theta) & \cos(\Theta) \end{bmatrix}.$$

The proof is in Appendix C.1.

Proposition 2 shows that the discretized complex SSM has an equivalent real SSM with doubled state dimension ( $N$ ), and a block-diagonal transition matrix multiplied with a scalar decay, where each  $2 \times 2$  block is a data-dependent rotation matrix ( $e_t^{\Delta_t A} \mathbf{R}_t$ ). We now show that the rotations can equivalently be absorbed into the input and output projections  $\mathbf{B}_t, \mathbf{C}_t$ , yielding an equivalent view that *complex SSMs are real SSMs equipped with data-dependent rotary embeddings (RoPE)*.

**Proposition 3** (Complex SSM, Data-Dependent RoPE Equivalence). *Under the notation established in Proposition 2, consider the real SSM defined in Eq. 7 unrolled for  $T$  time-steps. The output of the above SSM is equivalent to that of a vanilla scalar transition matrix-based SSM (Eq. 2) with a data-dependent rotary embedding applied on the  $\mathbf{B}, \mathbf{C}$  components of the SSM defined as:*

$$\mathbf{h}_t = e^{\Delta_t A_t} \mathbf{h}_{t-1} + \left(\prod_{i=0}^t \mathbf{R}_i^\top\right) \mathbf{B}_t x_t, \quad y_t = \left(\left(\prod_{i=0}^t \mathbf{R}_i^\top\right) \mathbf{C}_t\right)^\top \mathbf{h}_t \quad (8)$$

where the matrix production represents right matrix multiplication, e.g.,  $\prod_{i=0}^1 \mathbf{R}_i = \mathbf{R}_0 \mathbf{R}_1$ . We denote employing the vanilla SSM to compute the Complex SSM as “RoPE trick”.

The proof is in Appendix C.2.

To observe the connection of complex SSMs to RoPE embeddings, note that in the above proposition, the data-dependent rotations  $\mathbf{R}_i$  are aggregated across time-steps and applied to  $\mathbf{C}, \mathbf{B}$ , which, by the State Space Duality of Dao & Gu (2024), correspond to the Query ( $\mathbf{Q}$ ) and Key ( $\mathbf{K}$ ) components of Attention. Analogously, vanilla RoPE (Su et al., 2023) applies *data-independent* rotation matrices, where the rotation angles follow a fixed frequency schedule  $\boldsymbol{\theta}[i] = 10000^{-2i/N}$ .

**Remark 3** (Generality). Proposition 3 extends to the fully general case where the transition is given by any complex matrix. By the complex diagonalization theorem, such a matrix is unitarily equivalent to a complex diagonal matrix,  $\text{Diag}(\mathbf{A}(t) + i\boldsymbol{\theta}(t))$  with  $\mathbf{A}(t) \in \mathbb{R}^N$ . However, in practice, we restrict  $\mathbf{A}(t)$  to a scalar, mirroring the simplification from Mamba to Mamba-2, to enable faster implementation by avoiding GPU memory bottlenecks.

**Proposition 4** (Rotary Embedding Equivalence with Trapezoidal Discretization). *Discretizing a complex SSM with the trapezoidal rule (Proposition 1) yields the recurrence*

$$\begin{aligned} \mathbf{h}_t &= \alpha_t \mathbf{h}_{t-1} + \beta_t \left(\prod_{i=0}^{t-1} \mathbf{R}_i^\top\right) \mathbf{B}_{t-1} x_{t-1} + \gamma_t \left(\prod_{i=0}^t \mathbf{R}_i^\top\right) \mathbf{B}_t x_t, \\ y_t &= \left(\left(\prod_{i=0}^t \mathbf{R}_i^\top\right) \mathbf{C}_t\right)^\top \mathbf{h}_t. \end{aligned} \quad (9)$$

Here  $\mathbf{R}_t$  is the block-diagonal rotation matrix defined in Proposition 3.

The proof is in Appendix C.5.

**Remark 4** (RoPE Trick). Complex SSMs discretized with the general trapezoidal rule of a complex SSM naturally admit the RoPE trick we established for SSMs discretized with Euler’s rule.

### 3.3 MULTI-INPUT, MULTI-OUTPUT

During the decoding phase of autoregressive inference, outputs are generated one token at a time, and performance is typically measured using in *Tokens generated Per Second (TPS)*. In this metric, sub-quadratic models, such as Mamba-2 (Dao & Gu, 2024), have a significant advantage over standard Transformer-style attention, since they feature a fixed-size hidden state (Equation (2)) rather than maintaining a key-value (KV) cache that grows linearly with the sequence length.

TPS, however, does not explicitly factor in hardware efficiency, where we aim to be in a compute-bound regime (as opposed to memory-bound) in order to fully utilize on-chip accelerators. To better characterize hardware efficiency, we would need to consider the arithmetic intensity of token generation. Recall that arithmetic intensity is defined as FLOPs divided by the number of input-output bytes, for a given op. In order to fully utilize both the accelerators and the bandwidth, we would like the arithmetic intensity to match the ops:byte ratio of the hardware, which in the case of NVIDIA H100-SXM5, is 295.2 bfloat16 ops per second with respect to the DRAM, and 31.9 bfloat16 ops per second with respect to the SRAM [Fleetwood].

Table 2(a) shows the arithmetic intensity for a single generation in the SSM component of Mamba (with respect to 2-byte data). We see that it falls far short of a compute-bound regime, and moreover it is not clear how one can adjust the existing parameters in Mamba to mitigate the lack of hardware efficiency. We note that this observation applies generally to other sub-quadratic models, such as causal linear attention.

Input	Output	FLOPs	Arithmetic Intensity	Input	Output	FLOPs	Arithmetic Intensity
$H_t : (n, p)$	$y_t : (p)$	$5pn$	$\frac{5pn}{2(1+2n+p+np)}$	$H_t : (n, p)$	$y_t : (p, r)$	$4nrp + 2np$	$\frac{p(4nr+2n)}{2(1+2nr+pr+np)}$
$x_t : (p)$			$\approx 2.5 = \Theta(1)$	$x_t : (p, r)$			$\approx 2r = \Theta(r)$
$a_t : (1)$				$a_t : (1)$			
$b_t : (n)$				$b_t : (n, r)$			
$c_t : (n)$				$c_t : (n, r)$			

(a) SISO (2-byte data).

(b) MIMO (2-byte data).

Figure 2: Arithmetic Intensity for (a) SISO, (b) MIMO. Batch and head dimensions cancel out.

In light of this, we made the following simple adjustment to our recurrent relation: instead of transforming the input  $\mathbf{x}_t \in \mathbb{R}^p$  to state  $\mathbf{H}_t \in \mathbb{R}^{n \times p}$  via an outer product, i.e.,  $\mathbf{H}_t \leftarrow a_t \mathbf{H}_{t-1} + \mathbf{b}_t \otimes \mathbf{x}_t$ , we made such a transformation via a matrix product, i.e.,  $\mathbf{H}_t \leftarrow a_t \mathbf{H}_{t-1} + \mathbf{B}_t \mathbf{X}_t^\top$ , where  $\mathbf{B}_t \in \mathbb{R}^{n \times r}$  and  $\mathbf{X}_t \in \mathbb{R}^{p \times r}$  are now matrices with an addition rank  $r$ . The emission from state to output similarly acquire an extra rank  $r$ , i.e.,  $\mathbf{Y}_t \in \mathbb{R}^{r \times p} \leftarrow \mathbf{C}_t^\top \mathbf{H}_t$ , where  $\mathbf{C}_t \in \mathbb{R}^{n \times r}$ ,  $\mathbf{H}_t \in \mathbb{R}^{n \times p}$ . This simple change increases the arithmetic intensity of recurrence, which now scales with the rank  $r$  (Figure 2(b)). Hence, by increasing  $r$ , arithmetic intensity improves and shifts decode generation towards a more compute-bound regime. This increase in FLOPs during decode does not compromise runtime, as the operation is bounded by the I/O of state  $\mathbf{H}_t \in \mathbb{R}^{n \times p}$ .

Moreover, moving from outer-product-based state update to matrix-product-based coincides exactly with generalizing from SISO to MIMO SSM, with the rank  $r$  being the MIMO rank. Such a generalization recovers a key expressive feature of SSMs in classical literature; indeed, there has been previous work, namely Smith et al. (2023), that explored MIMO SSM as a drop-in replacement of attention, albeit not in the context of Mamba and not necessarily with inference in view.

Details of the MIMO formulation for Mamba-3 are provided in Appendix D.

### 3.4 MAMBA-3 ARCHITECTURE

The Mamba-3 block retains the overall layout of its predecessor while introducing several key modifications. Most notably, the SSD layer is replaced with the more expressive trapezoidal SSM defined in Proposition 4. The extra normalization layer, first introduced between Mamba-1 and Mamba-2 for training stability, is repositioned to follow the  $\mathbf{B}, \mathbf{C}$  projection, mirroring the QK-Norm commonly used in modern Transformers (Henry et al., 2020; Wortsman et al., 2023). Building on the findings of Yu & Erichson (2025), which prove adding channel-specific bias to  $\mathbf{B}$  in a blockwise variant of Mamba-1 grants universal approximation capabilities, Mamba-3 incorporates a head-specific, channel-wise bias into both the  $\mathbf{B}$  and  $\mathbf{C}$  components after its normalization. Our trapezoidal discretization complements this bias, eliminating the need for the original short causal convolution and its accompanying activation function (Section 4.3). Mamba-3 employs the SISO SSM by default, though we view its MIMO variant as a flexible option that can be toggled

Table 1: Downstream language modeling evaluations on models trained with 100B FineWeb-Edu tokens. Best results for each size are **bolded**, and second best are underlined. All models are trained with the same procedure. Mamba-3 outperforms Mamba-2 and others at every model scale.

Model	FW-Edu ppl ↓	LAMB. ppl ↓	LAMB. acc ↑	HellaS. acc_n ↑	PIQA acc ↑	Arc-E acc ↑	Arc-C acc_n ↑	WinoGr. acc ↑	OBQA acc_n ↑	Average acc ↑
Transformer-180M	16.89	45.0	32.5	39.0	<b>67.1</b>	59.8	27.9	51.2	21.8	42.8
Gated DeltaNet-180M	<u>16.61</u>	<b>35.9</b>	<b>33.7</b>	<u>40.2</u>	66.8	59.6	<b>28.5</b>	51.2	21.6	<u>43.1</u>
Mamba-2-180M	<u>16.76</u>	41.8	30.9	40.1	<u>66.8</u>	<u>60.1</u>	27.3	<b>52.0</b>	<b>23.2</b>	42.9
<b>Mamba-3-180M</b>	<b>16.59</b>	<u>37.7</u>	<u>32.5</u>	<b>40.8</b>	66.1	<b>61.5</b>	<u>27.9</u>	<b>52.0</b>	<u>22.8</u>	<b>43.4</b>
Transformer-440M	13.03	21.2	<b>41.7</b>	50.5	69.9	67.6	<u>34.6</u>	<b>56.7</b>	<b>26.0</b>	49.6
Gated DeltaNet-440M	13.12	<b>19.0</b>	40.4	50.5	70.5	67.5	34.0	55.3	25.8	<u>49.1</u>
Mamba-2-440M	13.00	19.6	40.8	<b>51.7</b>	70.6	68.8	<b>35.0</b>	54.1	<b>26.0</b>	49.6
<b>Mamba-3-440M</b>	<b>12.87</b>	<u>19.6</u>	<u>40.2</u>	<b>51.7</b>	<b>71.9</b>	<b>68.9</b>	34.4	<u>55.8</u>	<b>26.0</b>	<b>49.8</b>
Transformer-820M	11.42	15.0	44.7	57.2	72.6	71.6	39.2	57.7	26.8	52.8
Gated DeltaNet-820M	11.39	<b>12.7</b>	47.1	57.5	<u>72.6</u>	<u>72.5</u>	38.8	<u>57.9</u>	<b>30.6</b>	53.9
Mamba-2-820M	<u>11.35</u>	13.8	<u>45.0</u>	<u>58.1</u>	<u>72.5</u>	<u>72.3</u>	38.7	<u>56.8</u>	<u>30.2</u>	<u>53.4</u>
<b>Mamba-3-820M</b>	<b>11.23</b>	<u>12.9</u>	<b>47.2</b>	<b>58.8</b>	<b>73.6</b>	<b>72.7</b>	<b>40.2</b>	<b>58.4</b>	30.0	<b>54.4</b>
Transformer-1.5B	10.51	11.1	<b>50.3</b>	60.6	73.8	74.0	40.4	58.7	29.6	55.4
Gated DeltaNet-1.5B	10.51	<b>10.8</b>	49.9	60.5	<b>74.3</b>	73.3	40.4	<b>61.5</b>	30.4	<u>55.7</u>
Mamba-2-1.5B	10.47	12.0	47.8	61.4	73.6	75.3	41.8	57.5	<b>32.6</b>	<u>55.7</u>
<b>Mamba-3-1.5B</b>	<b>10.35</b>	<u>10.9</u>	49.4	<b>61.9</b>	73.6	<b>75.9</b>	<b>42.7</b>	<u>59.4</u>	<u>32.0</u>	<b>56.4</b>

depending on inference requirements. The overall architecture follows the Llama design (Grattafiori et al., 2024), alternating Mamba-3 and SwiGLU blocks with pre-normalization.

## 4 EMPIRICAL VALIDATION

We empirically validate our SSM-centric methodological changes through the overall Mamba-3 model on a host of synthetic and real world tasks. Section 4.1 compares our SISO-variant of Mamba-3 on language modeling and retrieval-based tasks, while Section 4.2 demonstrates inference efficiency of Mamba-3, and MIMO Mamba-3’s benefits over SISO Mamba-3 under fixed inference compute. We ablate the impact of our new discretization and BC bias on performance and show that complexification of the SSM leads to previously out-of-reach capabilities in Section 4.3.

### 4.1 LANGUAGE MODELING

All models are pretrained with 100B tokens of the FineWeb-Edu dataset (Penedo et al., 2024) with the Llama-3.1 tokenizer (Grattafiori et al., 2024) at a 2K context length with the same standard training protocol. Training and evaluation details can be found in Appendix E.

Across all four model scales, Mamba-3 outperforms popular baselines at various downstream tasks (Table 1). We highlight that Mamba-3 does not utilize the short convolution that has been empirically identified as an important component in many performant linear models (Allen-Zhu, 2025).

#### 4.1.1 RETRIEVAL CAPABILITIES

Beyond standard language modeling, an important measure for linear models is their retrieval ability — how well they can recall information from earlier in the sequence (Arora et al., 2025a;b). Unlike attention models, which can freely revisit past context with the growing KV cache, linear models must compress context into a fixed-size state. This trade-off is reflected in the Transformer baseline’s substantially stronger retrieval scores. To evaluate Mamba-3 under this lens, Table 2 compares it against baselines on both real-world and synthetic needle-in-a-haystack (NIAH) tasks (Hsieh et al., 2024), using our pretrained 1.5B models from Section 4.1. We restrict the task sequence length to 2K tokens to match the training setup and adopt the cloze-style format for our real-world tasks to mirror the next-token-prediction objective, following Arora et al. (2025b; 2024).

Mamba-3 is competitive on real-world associative recall and question-answering but struggles when extracting information from semi-structured or unstructured data. On synthetic NIAH tasks, however, Mamba-3 surpasses or matches baselines on most cases and notably demonstrates markedly better out-of-distribution retrieval abilities than its Mamba-2 predecessor.

#### 4.2 INFERENCE EFFICIENCY

In this section, we investigate our methodological changes in the context of inference performance. We first present our inference benchmark in Section 4.2.1; we then establish a framework for comparing the inference performance in Section 4.2.2. Finally, we focus on the effectiveness of MIMO in Section 4.2.3.

Table 2: Retrieval capabilities measured by a mixture of real-world and synthetic retrieval tasks. Real-world retrieval tasks utilize cloze variants of the original datasets and are truncated to 2K length. Mamba-3 demonstrates strong associative recall and question-answering but suffers with information extraction of semi-structured and unstructured data. Mamba-3 has strong needle-in-a-haystack (NIAH) accuracy and generalizes outside its trained context.

Model (1.5B)	SWDE	SQUAD	FDA	TQA	NQ	Drop	NIAH-Single-1			NIAH-Single-2			NIAH-Single-3		
Context Length	2048						1024	2048	4096	1024	2048	4096	1024	2048	4096
Transformer	48.9	46.6	58.4	67.5	31.7	26.4	100.0	100.0	0.0	92.2	100.0	0.0	98.6	99.4	0
Gated DeltaNet	<b>32.7</b>	40.0	<b>28.3</b>	63.5	25.7	24.5	<b>100.0</b>	<b>100.0</b>	<b>99.8</b>	<b>100.0</b>	93.8	49.8	83.8	68.4	<b>34.2</b>
Mamba-2	30.7	39.1	23.7	64.3	25.1	<b>28.5</b>	<b>100.0</b>	99.6	62.0	<b>100.0</b>	53.8	11.8	<b>95.8</b>	<b>87.4</b>	13.4
<b>Mamba-3</b>	28.5	<b>40.1</b>	23.4	<b>64.5</b>	<b>26.5</b>	<u>27.4</u>	<b>100.0</b>	<b>100.0</b>	<u>88.2</u>	<b>100.0</b>	<b>95.4</b>	<b>50.6</b>	<u>92.4</u>	<u>81.4</u>	<b>34.2</b>

Model	FP32		BF16	
	$d_{\text{state}}=64$	$d_{\text{state}}=128$	$d_{\text{state}}=64$	$d_{\text{state}}=128$
Mamba-2	0.295	0.409	0.127	0.203
Gated DeltaNet	0.344	0.423	0.176	0.257
Mamba-3 (SISO)	0.261	0.356	0.106	0.152
Mamba-3 (MIMO)	0.285	0.392	0.136	0.185

Table 3: Latency (in milliseconds) comparison across models, precision, and  $d_{\text{state}}$  values. Both Mamba-3 SISO and MIMO are faster than the Mamba-2 and Gated DeltaNet at the commonly used bf16,  $d_{\text{state}}=128$  setting.

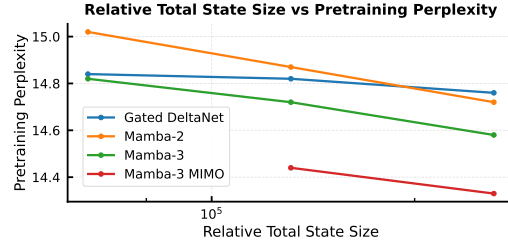


Figure 3: Exploration of state size (inference speed proxy) versus pretraining perplexity (performance proxy). Mamba-3 MIMO drives the-Pareto frontier without increasing state size.

#### 4.2.1 MAMBA-3 IS FAST AT INFERENCE

We benchmarked the wallclock time for a single decoding step, and in a single sequence mixing layer, for each subquadratic model we consider in this paper. We adopted the standard reference code for Mamba-2 and GDN, while writing our own custom kernels for the Mamba-3 step function. The result is recorded in Table 3. We observe that, despite having a more sophisticated SSM structure, Mamba-3 is in fact noticeably faster when compared to Mamba-2, which in turn is faster than GDN, illustrating the viability of our inference-first approach.

#### 4.2.2 A PARETO FRONT FOR INFERENCE EFFICIENCY

For Mamba and many variants of sub-quadratic models, the generation of tokens during decoding is heavily dominated by memory I/O due to the low arithmetic intensity of computing the recurrent update (c.f. Section 3.3). Furthermore, among the data being transferred, the latent state  $\mathbf{H}_t$  dominates in terms of size. Indeed, from Table 3, we see that the runtime scales with  $d_{\text{state}}$ , which configures the size of the hidden state.

As  $d_{\text{state}}$  dominates the decode runtime for the subquadratic models considered in this paper, we opt to use it as a proxy for inference speed. By plotting the validation perplexity (itself a proxy for model performance) as a function of  $d_{\text{state}}$ , we aim to formulate a holistic picture about how the subquadratic models can trade off performance with inference speed.

Figure 3<sup>1</sup> shows such a Pareto front for the subquadratic models considered in this paper. For each data point, we train a 440M parameter model to 17.8 billion tokens on the Fineweb-Edu dataset, where the model is configured with a  $d_{\text{state}}$  of  $\{32, 64, 128\}$ . As expected, we observe an inverse correlation between validation loss and  $d_{\text{state}}$ ; moreover, we noticed a general downward shift on the Pareto front moving from Mamba-2 to Mamba-3. A further downward shift is observed when moving from the SISO variant of Mamba-3 to the MIMO variant of Mamba-3 (where we set the MIMO rank  $r=4$  and decrease our MLP inner dimension to match the parameter count). This highlights both the expressivity gain coming our methodology change as well as the effectiveness of the MIMO mechanism in improving decoding efficiency.

<sup>1</sup>The  $d_{\text{state}}=32$  Mamba-3 MIMO did not finish training at time of submission.



**Table 4: Left:** Ablations on core modeling components of Mamba-3, results on test split of dataset. A combination of our BC bias and trapezoidal discretization makes the convolution optional. **Right:** Formal language evaluation (scaled accuracy, %). Higher is better. Models are trained on short sequences and evaluated on longer lengths to test length generalization. For Gated DeltaNet we report the variant with eigenvalue range  $[-1, 1]$ .

Model Variant	ppl ↓	Model	Parity ↑	Arith. w/o ↑ brackets	Arith. w/ ↑ brackets
Mamba-3 — bias — trap	16.68	Mamba-3	100.00	98.51	87.75
Mamba-3 — bias	16.49	Mamba-3 (w/o RoPE)	2.27	1.49	0.72
Mamba-3	<b>15.72</b>	Mamba-2	0.90	47.81	0.88
Mamba-3 + conv	15.85	Gated DeltaNet $[-1, 1]$	100.00	99.25	93.50

(a) Component ablation (350M).

(b) Performance comparison of various models on formal language tasks. Results show that unlike Mamba-2, Mamba-3 features state tracking ability stemming from data-dependent RoPE embeddings.

#### 4.2.3 MIMO ENHANCES INFERENCE EFFICIENCY

With higher arithmetic intensity, MIMO increases the decoding FLOPs. Moreover, from Table 3, we observe that the additional FLOPs does not have a drastic impact on the decode runtime.<sup>2</sup> The implication is that any performance gain from MIMO translates to efficiency gain in decoding, and this is in fact supported by the downward shift of the MIMO Pareto curve we observed in Section 4.2.2.

We aim to further verify the gain from MIMO by investigating its language-modeling capabilities. To that end, we train a 440M parameter MIMO model with MIMO rank  $r = 4$  on 100B tokens on Fineweb-Edu (i.e., same setting as the 440M parameter run in Section 4.1; we did not train 820M or 1.5B model due to compute constraints). To ensure the total parameter count equals SISO, we decrease the inner dimension of the MLP layers to compensate for the increase due to the MIMO projections.

On both validation perplexity and our suite of language evaluation tasks (Table 5), we see significant gain when moving from SISO to MIMO. Namely, we attain a perplexity gain of 0.16 on the 100B tokens run, and Figure 3 illustrates the downward shift in our validation loss. On the language evaluation front, we see significant gain on most tasks when compared to SISO, resulting in an overall gain of 1.2 point over SISO. This strongly supports MIMO as a SSM-centric technique to improve model quality without compromising decoding speed.

#### 4.3 SSM-CENTRIC METHODOLOGICAL ABLATIONS

Table 4a ablates the changes made to the core SSM component, mainly the introduction of BC bias and trapezoidal discretization. We report the pretraining test perplexity on models at the 440M scale, trained for Chinchilla optimal tokens. We find that the bias and trapezoidal SSM synergize well and make the short convolution utilized by many current linear models redundant.

We empirically demonstrate that data-dependent RoPE in Mamba-3 enables state tracking. Following Grazi et al. (2025), we evaluate on tasks from the Chomsky hierarchy—Parity, Modular Arithmetic (without brackets), and Modular Arithmetic (with brackets)—and report scaled accuracies in Table 4b. Mamba-3 solves Parity and Modular Arithmetic (without brackets), and nearly closes the accuracy gap on Modular Arithmetic (with brackets). In contrast, Mamba-3 without RoPE and Mamba-2 fail to learn these tasks. We use the state-tracking-enabled *Gated DeltaNet* variant of Grazi et al. (2025) and observe that *Mamba-3* is competitive—matching parity and approaching its performance on both modular-arithmetic tasks. Experimental settings are covered in Appendix E.

## 5 CONCLUSION AND FUTURE WORK

We introduce Mamba-3, an SSM model with three axes of improvement rooted in SSM principles: (i) *improved quality*, via trapezoidal discretization; (ii) *new capabilities*, through complex SSMs that recover state-tracking; and (iii) *higher inference efficiency*, with a MIMO formulation that raises arithmetic intensity. Mamba-3 delivers strong language modeling results and establishes a new Pareto frontier on the performance-efficiency axes with respect to strong baseline models. A limitation remains in retrieval, where fixed-state architectures lags attention-based models. We see **hybrid Mamba-3 architectures** that integrate retrieval mechanisms as a promising path, alongside broader application of our design principles to linear-time sequence models.

<sup>2</sup>The kernel for MIMO Mamba-3 in fact fuses the MIMO projection, and so the reported wallclock time is actually an overestimate for the pure SSM update.

## REFERENCES

- Zeyuan Allen-Zhu. Physics of Language Models: Part 4.1, Architecture Design and the Magic of Canon Layers. *SSRN Electronic Journal*, May 2025. <https://ssrn.com/abstract=5240330>.
- Aryaman Arora, Neil Rathi, Nikil Roashan Selvam, Róbert Csordás, Dan Jurafsky, and Christopher Potts. Mechanistic evaluation of transformers and state space models, 2025a. URL <https://arxiv.org/abs/2505.15105>.
- Simran Arora, Aman Timalsina, Aaryan Singhal, Benjamin Spector, Sabri Eyuboglu, Xinyi Zhao, Ashish Rao, Atri Rudra, and Christopher Ré. Just read twice: closing the recall gap for recurrent language models, 2024. URL <https://arxiv.org/abs/2407.05483>.
- 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, 2025b. URL <https://arxiv.org/abs/2402.18668>.
- Aviv Bick, Kevin Y. Li, Eric P. Xing, J. Zico Kolter, and Albert Gu. Transformers to ssms: Distilling quadratic knowledge to subquadratic models, 2025a. URL <https://arxiv.org/abs/2408.10189>.
- Aviv Bick, Eric Xing, and Albert Gu. Understanding the skill gap in recurrent language models: The role of the gather-and-aggregate mechanism, 2025b. URL <https://arxiv.org/abs/2504.18574>.
- Yonatan Bisk, Rowan Zellers, Ronan Le Bras, Jianfeng Gao, and Yejin Choi. Piqa: Reasoning about physical commonsense in natural language, 2019. URL <https://arxiv.org/abs/1911.11641>.
- Krzysztof Choromanski, Valerii Likhoshesterov, David Dohan, Xingyou Song, Andreea Gane, Tamas Sarlos, Peter Hawkins, Jared Davis, Afroz Mohiuddin, Lukasz Kaiser, David Belanger, Lucy Colwell, and Adrian Weller. Rethinking attention with performers, 2022. URL <https://arxiv.org/abs/2009.14794>.
- 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, 2024. URL <https://arxiv.org/abs/2405.21060>.
- Dheeru Dua, Yizhong Wang, Pradeep Dasigi, Gabriel Stanovsky, Sameer Singh, and Matt Gardner. Drop: A reading comprehension benchmark requiring discrete reasoning over paragraphs, 2019. URL <https://arxiv.org/abs/1903.00161>.
- Christopher Fleetwood. Domain specific architectures for ai inference. URL <https://fleetwood.dev/posts/domain-specific-architectures>.
- 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 Muennighoff, 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>.
- Madan Gopal. *Modern control system theory*. New Age International, 1993.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, and et. al. The llama 3 herd of models, 2024. URL <https://arxiv.org/abs/2407.21783>.
- Riccardo Grazi, Julien Siems, Simon Schrod, Thomas Brox, and Frank Hutter. Is mamba capable of in-context learning?, 2024. URL <https://arxiv.org/abs/2402.03170>.
- Riccardo Grazi, Julien Siems, Arber Zela, Jörg K. H. Franke, Frank Hutter, and Massimiliano Pontil. Unlocking state-tracking in linear rnns through negative eigenvalues, 2025. URL <https://arxiv.org/abs/2411.12537>.

- Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces, 2024. URL <https://arxiv.org/abs/2312.00752>.
- Albert Gu, Karan Goel, and Christopher Ré. Efficiently modeling long sequences with structured state spaces, 2022a. URL <https://arxiv.org/abs/2111.00396>.
- Albert Gu, Ankit Gupta, Karan Goel, and Christopher Ré. On the parameterization and initialization of diagonal state space models. *arXiv preprint arXiv:2206.11893*, 2022b. URL <https://arxiv.org/abs/2206.11893>.
- Ankit Gupta, Albert Gu, and Jonathan Berant. Diagonal state spaces are as effective as structured state spaces, 2022. URL <https://arxiv.org/abs/2203.14343>.
- Alex Henry, Prudhvi Raj Dachapally, Shubham Pawar, and Yuxuan Chen. Query-key normalization for transformers, 2020. URL <https://arxiv.org/abs/2010.04245>.
- Cheng-Ping Hsieh, Simeng Sun, Samuel Kriman, Shantanu Acharya, Dima Rekesh, Fei Jia, Yang Zhang, and Boris Ginsburg. Ruler: What’s the real context size of your long-context language models?, 2024. URL <https://arxiv.org/abs/2404.06654>.
- Samy Jelassi, David Brandfonbrener, Sham M. Kakade, and Eran Malach. Repeat after me: Transformers are better than state space models at copying, 2024. URL <https://arxiv.org/abs/2402.01032>.
- Mandar Joshi, Eunsol Choi, Daniel S. Weld, and Luke Zettlemoyer. Triviaqa: A large scale distantly supervised challenge dataset for reading comprehension, 2017. URL <https://arxiv.org/abs/1705.03551>.
- Rudolph Emil Kalman. A new approach to linear filtering and prediction problems. 1960.
- Angelos Katharopoulos, Apoorv Vyas, Nikolaos Pappas, and François Fleuret. Transformers are rnns: Fast autoregressive transformers with linear attention, 2020. URL <https://arxiv.org/abs/2006.16236>.
- Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, Kristina Toutanova, Llion Jones, Matthew Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob Uszkoreit, Quoc Le, and Slav Petrov. Natural questions: A benchmark for question answering research. *Transactions of the Association for Computational Linguistics*, 7:452–466, 2019. doi: 10.1162/tacl.a.00276. URL <https://aclanthology.org/Q19-1026/>.
- Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model serving with pagedattention, 2023. URL <https://arxiv.org/abs/2309.06180>.
- Baolin Li, Yankai Jiang, Vijay Gadepally, and Devesh Tiwari. Llm inference serving: Survey of recent advances and opportunities, 2024. URL <https://arxiv.org/abs/2407.12391>.
- William Merrill, Jackson Petty, and Ashish Sabharwal. The illusion of state in state-space models, 2025. URL <https://arxiv.org/abs/2404.08819>.
- Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. Can a suit of armor conduct electricity? a new dataset for open book question answering, 2018. URL <https://arxiv.org/abs/1809.02789>.
- Team OLMo, Pete Walsh, Luca Soldaini, Dirk Groeneveld, Kyle Lo, Shane Arora, Akshita Bhagia, Yuling Gu, Shengyi Huang, Matt Jordan, Nathan Lambert, Dustin Schwenk, Oyvind Tafjord, Taira Anderson, David Atkinson, Faeze Brahman, Christopher Clark, Pradeep Dasigi, Nouha Dziri, Michal Guerquin, and et. al. 2 olmo 2 furious, 2025. URL <https://arxiv.org/abs/2501.00656>.
- Antonio Orvieto, Samuel L Smith, Albert Gu, Anushan Fernando, Caglar Gulcehre, Razvan Pascanu, and Soham De. Resurrecting recurrent neural networks for long sequences, 2023. URL <https://arxiv.org/abs/2303.06349>.
- Daniele Paliotta, Junxiong Wang, Matteo Pagliardini, Kevin Y. Li, Aviv Bick, J. Zico Kolter, Albert Gu, François Fleuret, and Tri Dao. Thinking slow, fast: Scaling inference compute with distilled reasoners, 2025. URL <https://arxiv.org/abs/2502.20339>.

- Denis Paperno, Germán Kruszewski, Angeliki Lazaridou, Quan Ngoc Pham, Raffaella Bernardi, Sandro Pezzelle, Marco Baroni, Gemma Boleda, and Raquel Fernández. The lambada dataset: Word prediction requiring a broad discourse context, 2016. URL <https://arxiv.org/abs/1606.06031>.
- Jongho Park, Jaeseung Park, Zheyang Xiong, Nayoung Lee, Jaewoong Cho, Samet Oymak, Kangwook Lee, and Dimitris Papailiopoulos. Can mamba learn how to learn? a comparative study on in-context learning tasks, 2024. URL <https://arxiv.org/abs/2402.04248>.
- Guilherme Penedo, Hynek Kydlíček, Loubna Ben allal, Anton Lozhkov, Margaret Mitchell, Colin Raffel, Leandro Von Werra, and Thomas Wolf. The fineweb datasets: Decanting the web for the finest text data at scale, 2024. URL <https://arxiv.org/abs/2406.17557>.
- Bo Peng, Ruichong Zhang, Daniel Goldstein, Eric Alcaide, Xingjian Du, Haowen Hou, Jiaju Lin, Jiaxing Liu, Janna Lu, William Merrill, Guangyu Song, Kaifeng Tan, Saiteja Utpala, Nathan Wilce, Johan S. Wind, Tianyi Wu, Daniel Wuttke, and Christian Zhou-Zheng. Rwkv-7 ”goose” with expressive dynamic state evolution, 2025. URL <https://arxiv.org/abs/2503.14456>.
- Pranav Rajpurkar, Jian Zhang, and Percy Liang. Know what you don’t know: Unanswerable questions for squad. In *ACL 2018*, 2018.
- Yuval Ran-Milo, Eden Lumbroso, Edo Cohen-Karlik, Raja Giryes, Amir Globerson, and Nadav Cohen. Provable benefits of complex parameterizations for structured state space models, 2024. URL <https://arxiv.org/abs/2410.14067>.
- 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>.
- Yash Sarrof, Yana Veitsman, and Michael Hahn. The expressive capacity of state space models: A formal language perspective, 2024. URL <https://arxiv.org/abs/2405.17394>.
- Imanol Schlag, Kazuki Irie, and Jürgen Schmidhuber. Linear transformers are secretly fast weight programmers, 2021. URL <https://arxiv.org/abs/2102.11174>.
- Julien Siems, Timur Carstensen, Arber Zela, Frank Hutter, Massimiliano Pontil, and Riccardo Grazi. Deltaproduct: Improving state-tracking in linear rnns via householder products, 2025. URL <https://arxiv.org/abs/2502.10297>.
- Jimmy T. H. Smith, Andrew Warrington, and Scott W. Linderman. Simplified state space layers for sequence modeling, 2023. URL <https://arxiv.org/abs/2208.04933>.
- Charlie Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. Scaling llm test-time compute optimally can be more effective than scaling model parameters, 2024. URL <https://arxiv.org/abs/2408.03314>.
- Jianlin Su, Yu Lu, Shengfeng Pan, Ahmed Murtadha, Bo Wen, and Yunfeng Liu. Roformer: Enhanced transformer with rotary position embedding, 2023. URL <https://arxiv.org/abs/2104.09864>.
- 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, 2023. URL <https://arxiv.org/abs/2307.08621>.
- Endre Süli and David F. Mayers. *An Introduction to Numerical Analysis*. Cambridge University Press, 2003.
- Gemma Team, Aishwarya Kamath, Johan Ferret, Shreya Pathak, Nino Vieillard, Ramona Merhej, Sarah Perrin, Tatiana Matejovicova, Alexandre Ramé, Morgane Rivière, Louis Rouillard, Thomas Mesnard, Geoffrey Cideron, Jean bastien Grill, Sabela Ramos, Edouard Yvinec, Michelle Casbon, Etienne Pot, Ivo Penchev, Gaël Liu, and et. al. Gemma 3 technical report, 2025. URL <https://arxiv.org/abs/2503.19786>.
- M. Tenenbaum and H. Pollard. *Ordinary Differential Equations: An Elementary Textbook for Students of Mathematics, Engineering, and the Sciences*. Dover Books on Mathematics. Dover Publications, 1985. ISBN 9780486649405. URL <https://books.google.com/books?id=iU4zDAAAQBAJ>.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pp. 5998–6008, 2017. URL <http://arxiv.org/abs/1706.03762>.

- Johannes von Oswald, Nino Scherrer, Seijin Kobayashi, Luca Versari, Songlin Yang, Maximilian Schlegel, Kaitlin Maile, Yanick Schimpf, Oliver Sieberling, Alexander Meulemans, Rif A. Saurous, Guillaume Lajoie, Charlotte Frenkel, Razvan Pascanu, Blaise Agüera y Arcas, and João Sacramento. Mesanet: Sequence modeling by locally optimal test-time training, 2025. URL <https://arxiv.org/abs/2506.05233>.
- Mitchell Wortsman, Peter J. Liu, Lechao Xiao, Katie Everett, Alex Alemi, Ben Adlam, John D. Co-Reyes, Izzeddin Gur, Abhishek Kumar, Roman Novak, Jeffrey Pennington, Jascha Sohl-dickstein, Kelvin Xu, Jaehoon Lee, Justin Gilmer, and Simon Kornblith. Small-scale proxies for large-scale transformer training instabilities, 2023. URL <https://arxiv.org/abs/2309.14322>.
- Yangzhen Wu, Zhiqing Sun, Shanda Li, Sean Welleck, and Yiming Yang. Inference scaling laws: An empirical analysis of compute-optimal inference for problem-solving with language models, 2025. URL <https://arxiv.org/abs/2408.00724>.
- Songlin Yang, Jan Kautz, and Ali Hatamizadeh. Gated delta networks: Improving mamba2 with delta rule, 2025a. URL <https://arxiv.org/abs/2412.06464>.
- Songlin Yang, Bailin Wang, Yu Zhang, Yikang Shen, and Yoon Kim. Parallelizing linear transformers with the delta rule over sequence length, 2025b. URL <https://arxiv.org/abs/2406.06484>.
- Annan Yu and N. Benjamin Erichson. Block-biased mamba for long-range sequence processing, 2025. URL <https://arxiv.org/abs/2505.09022>.
- 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>.

**LLM Usage.** We utilized Large Language Models to polish the writing in our submission as well as generate latex code for formatting tables and figures.

## A RELATED WORK

**Linear-time sequence mixers.** State-space models (SSMs) provide linear-time sequence mixing through explicit dynamical states and efficient scan/convolution implementations, offering significant computational advantages over quadratic-time attention mechanisms (Gu et al., 2022a; Smith et al., 2023; Gupta et al., 2022). Mamba-1 (Gu & Dao, 2024) introduced input-dependent selectivity to SSMs, while Mamba-2 (Dao & Gu, 2024) formalized the connection between SSMs and attention via structured state-space duality (SSD) (Katharopoulos et al., 2020; Choromanski et al., 2022). Despite matching transformers on standard language understanding benchmarks, these recurrent models exhibit limitations on tasks requiring precise algorithmic reasoning. Recent evaluations identified gaps in capabilities such as associative retrieval (Bick et al., 2025b; Arora et al., 2025a), exact copying (Jelassi et al., 2024), and in-context learning (Park et al., 2024; Grazzi et al., 2024). To address these limitations, DeltaNet enhances linear attention by replacing additive updates with delta-rule recurrence (Schlag et al., 2021), with recent work developing hardware-efficient, sequence-parallel training algorithms for this architecture (Yang et al., 2025b). This has catalyzed a broader effort to improve the algorithmic capabilities of linear-time models through architectural innovations including gating mechanisms, improved state transition dynamics, and hybrid approaches (Peng et al., 2025; Siems et al., 2025; Yang et al., 2025a; Paliotta et al., 2025; Bick et al., 2025a).

**Expressivity and state tracking in recurrent mixers.** Recent work characterizes the types of state that recurrent, constant-memory mixers can maintain, revealing algorithmic deficiencies in previous SSM-based models. Merrill et al. (2025) show that under finite precision, practical SSMs collapse to  $TC^0$ , leading to failures on tasks like permutation composition over  $S_5$  unless the primitive is extended. Similarly, Yu & Erichson (2025) prove that a single-layer Mamba is not a universal approximator. Several modifications have been proposed to improve expressivity. For instance, the same work shows that a block-biased variant regains the universal approximation property with only minor changes, either through block decomposition or a channel-specific bias. Allowing negative eigenvalues or non-triangular transitions enables linear RNNs—including diagonal and Householder/DeltaNet forms—to capture parity and, under mild assumptions, regular languages (Grazzi et al., 2025). Complex-valued parameterizations provide another avenue for enhanced expressivity. Diagonal LTI SSMs demonstrate effectiveness for language modeling (Gu et al., 2022b; Orvieto et al., 2023), with complex variants achieving equivalent functions using smaller, well-conditioned parameters (Ran-Milo et al., 2024). However, the introduction of selectivity—the central innovation of modern SSMs (Gu & Dao, 2024)—narrowed the performance gap with Transformers by enabling input-dependent dynamics and achieving state-of-the-art results on language modeling benchmarks, leading practitioners to abandon complex states in favor of simpler real-valued architectures. We extend this line of work by reintroducing complex-valued state evolution that yields a real SSM with doubled dimensionality and block-diagonal rotations applied to the update rule—analogue through SSD (Dao & Gu, 2024) to how RoPE (Su et al., 2023) applies complex rotations to queries and keys in attention. The resulting data-dependent rotational structure expands stable dynamics to include oscillatory modes, enabling richer states while maintaining constant memory and linear-time complexity.

## B TRAPEZOIDAL DISCRETIZATION PROOFS

### B.1 PROOF OF PROPOSITION 5

**Proposition 5** (Variation of Constants (Tenenbaum & Pollard, 1985)). *Consider the linear SSM,*

$$\dot{\mathbf{h}}(t) = A(t)\mathbf{h}(t) + \mathbf{B}(t)x(t),$$

where  $\mathbf{h}(t) \in \mathbb{R}^N$  is the hidden state,  $A(t) \in \mathbb{R}$  is the scalar state transition matrix and  $\mathbf{B}(t)x(t) \in \mathbb{R}^N$  is the input projection. For a time grid  $\tau_t = \tau_{t-1} + \Delta_t$ , the hidden state satisfies,

$$\mathbf{h}_t \approx e^{\Delta_t A_t} \mathbf{h}_{t-1} + \underbrace{\int_{\tau_{t-1}}^{\tau_t} e^{(\tau_t - \tau) A_t} \mathbf{B}(\tau) x(\tau) d\tau}_{\text{state-update}}, \quad (10)$$

*Proof.* Let  $\Phi(t, s)$  be the fundamental solution of the homogeneous system  $\dot{\mathbf{h}}(t) = A(t)\mathbf{h}(t)$ , i.e.,  $\partial_t \Phi(t, s) = A(t)\Phi(t, s)$  and  $\Phi(s, s) = 1$  (since  $A$  is scalar). By variation of constants,

$$\mathbf{h}(t) = \Phi(t, s)\mathbf{h}(s) + \int_s^t \Phi(t, \tau) \mathbf{B}(\tau) x(\tau) d\tau.$$

Choosing  $(s, t) = (t_{k-1}, t_k)$  gives the exact one-step relation

$$h_k = \Phi(t_k, t_{k-1})h_{k-1} + \int_{t_{k-1}}^{t_k} \Phi(t_k, \tau)B(\tau)x(\tau)d\tau.$$

On the step  $[t_{k-1}, t_k]$ , apply zero-order hold to  $A(\cdot)$ :  $A(\tau) \approx A_k$ . Then  $\Phi(t_k, \tau) \approx e^{(t_k - \tau)A_k}$  and  $\Phi(t_k, t_{k-1}) \approx e^{\Delta_k A_k}$ , yielding

$$h_k \approx e^{\Delta_k A_k} h_{k-1} + \int_{t_{k-1}}^{t_k} e^{(t_k - \tau)A_k} B(\tau)x(\tau)d\tau,$$

which is equation 10.  $\square$

## B.2 TRAPEZOID DISCRETIZATION' MASK MATRIX

*Proof.* When viewing the tensor contraction form, let us call  $C = (T, N), B = (S, N), L = (T, S), X = (S, P)$  based on the Mamba-2 paper. With this decomposition of our mask, we can view  $L = \text{contract}(TZ, ZS \rightarrow TS)(L_1, L_2)$ .

The original contraction can be seen as

$$\text{contract}(TN, SN, TS, SP \rightarrow TP)(C, B, L, X)$$

We can now view it as

$$\text{contract}(TN, SN, TJ, JS, SP \rightarrow TP)(C, B, L_1, L_2, X)$$

This can be broken into the following:

$$Z = \text{contract}(SN, SP \rightarrow SNP)(B, X)$$

$$Z' = \text{contract}(JS, SNP \rightarrow JNP)(L_2, Z)$$

$$H = \text{contract}(TJ, JNP \rightarrow TNP)(L_1, Z')$$

$$Y = \text{contract}(TN, TNP \rightarrow TP)(C, H)$$

Thus, we can view this step:  $\text{contract}(ZS, SNP \rightarrow ZNP)(L_2, Z)$  as a conv of size two applied on Bx with the traditional SSD  $L = L_1$  matrix.  $\square$

## C COMPLEX SSM PROOFS

### C.1 PROOF OF PROPOSITION 2

**Proposition 2** (Complex-to-Real SSM Equivalence). *Consider a complex-valued SSM*

$$\dot{\mathbf{h}}(t) = \text{Diag}(A(t) + i\theta(t))\mathbf{h}(t) + (\mathbf{B}(t) + i\hat{\mathbf{B}}(t))x(t), \quad (6)$$

$$y(t) = \text{Re}\left((\mathbf{C}(t) + i\hat{\mathbf{C}}(t))^\top \mathbf{h}(t)\right),$$

where  $\mathbf{h}(t) \in \mathbb{C}^{N/2}$ ,  $\theta(t), \mathbf{B}(t), \hat{\mathbf{B}}(t), \mathbf{C}(t), \hat{\mathbf{C}}(t) \in \mathbb{R}^{N/2}$ , and  $x(t), A(t) \in \mathbb{R}$ . Under Euler discretization, this system is equivalent to a real-valued SSM

$$\mathbf{h}_t = e^{\Delta_t A_t} \mathbf{R}_t \mathbf{h}_{t-1} + \Delta_t \mathbf{B}_t x_t, \quad (7)$$

$$y_t = \mathbf{C}_t^\top \mathbf{h}_t,$$

with state  $\mathbf{h}_t \in \mathbb{R}^N$ , projections

$$\mathbf{B}_t = \begin{bmatrix} \mathbf{B}_t \\ \hat{\mathbf{B}}_t \end{bmatrix} \in \mathbb{R}^N, \quad \mathbf{C}_t = \begin{bmatrix} \mathbf{C}_t \\ -\hat{\mathbf{C}}_t \end{bmatrix} \in \mathbb{R}^N,$$

and a transition matrix

$$\mathbf{R}_t = \text{Block}\left(\{R(\Delta_t \theta_t[i])\}_{i=1}^{N/2}\right) \in \mathbb{R}^{N \times N}, \quad R(\Theta) = \begin{bmatrix} \cos(\Theta) & -\sin(\Theta) \\ \sin(\Theta) & \cos(\Theta) \end{bmatrix}.$$

*Proof.* We first present the derivation for  $N = 2$ ; the block-diagonal structure for general even  $N$  follows by grouping pairs of coordinates.

Let  $h_t + i\hat{h}_t$  denote the complexified hidden state, with parameters  $A(t) + i\theta(t)$  and  $B(t) + i\hat{B}(t)$  for the transition and input, respectively. By the variation of constants formula (Proposition 5), applying zero-order hold and Euler's rule over a step  $[t_{k-1}, t_k]$  gives

$$h_k + i\hat{h}_k = e^{\Delta_t(A_t + i\theta_t)}(h_{k-1} + i\hat{h}_{k-1}) + \Delta_t(B_t + i\hat{B}_t)x_t.$$

Expanding the exponential,

$$e^{\Delta_t(A_t + i\theta_t)} = e^{\Delta_t A_t} \left( \cos(\Delta_t \theta_t) + i \sin(\Delta_t \theta_t) \right),$$

so in real coordinates  $\mathbf{h}_t = \begin{bmatrix} h_t \\ \hat{h}_t \end{bmatrix} \in \mathbb{R}^2$  the recurrence becomes

$$\mathbf{h}_t = e^{\Delta_t A_t} \underbrace{\begin{bmatrix} \cos(\Delta_t \theta_t) & -\sin(\Delta_t \theta_t) \\ \sin(\Delta_t \theta_t) & \cos(\Delta_t \theta_t) \end{bmatrix}}_{R(\Delta_t \theta_t)} \mathbf{h}_{t-1} + \Delta_t \begin{bmatrix} B_t \\ \hat{B}_t \end{bmatrix} x_t.$$

Stacking across  $N/2$  such pairs yields the block-diagonal transition

$$\mathbf{h}_t = e^{\Delta_t A_t} \text{Block}(\{R(\Delta_t \theta_t[i])\}_{i=1}^{N/2}) \mathbf{h}_{t-1} + \Delta_t \begin{bmatrix} \mathbf{B}_t \\ \hat{\mathbf{B}}_t \end{bmatrix} x_t.$$

For the output,

$$y_t = \text{Re}\left((\mathbf{C}_t + i\hat{\mathbf{C}}_t)^\top (h_t + i\hat{h}_t)\right) = \begin{bmatrix} \mathbf{C}_t \\ -\hat{\mathbf{C}}_t \end{bmatrix}^\top \mathbf{h}_t,$$

which defines the real projection  $\mathbf{C}_t \in \mathbb{R}^N$  in the proposition. This proves the equivalence between complex SSM and the real block-diagonal system with rotations.  $\square$

## C.2 PROOF OF PROPOSITION 3

**Proposition 3** (Complex SSM, Data-Dependent RoPE Equivalence). *Under the notation established in Proposition 2, consider the real SSM defined in Eq. 7 unrolled for  $T$  time-steps. The output of the above SSM is equivalent to that of a vanilla scalar transition matrix-based SSM (Eq. 2) with a data-dependent rotary embedding applied on the  $\mathbf{B}, \mathbf{C}$  components of the SSM defined as:*

$$\mathbf{h}_t = e^{\Delta_t A_t} \mathbf{h}_{t-1} + \left(\prod_{i=0}^t \mathbf{R}_i^\top\right) \mathbf{B}_t x_t, \quad \mathbf{y}_t = \left(\prod_{i=0}^t \mathbf{R}_i^\top\right) \mathbf{C}_t^\top \mathbf{h}_t \quad (8)$$

where the matrix production represents right matrix multiplication, e.g.,  $\prod_{i=0}^1 \mathbf{R}_i = \mathbf{R}_0 \mathbf{R}_1$ . We denote employing the vanilla SSM to compute the Complex SSM as “RoPE trick”.

*Proof.* Consider the SSM

$$\mathbf{h}_t = e^{\Delta_t A_t} \mathbf{R}_t \mathbf{h}_{t-1} + \mathbf{B}_t x_t, \quad \mathbf{y}_t = \mathbf{C}_t^\top \mathbf{h}_t, \quad (11)$$

where (as in Proposition 3)  $A_t \in \mathbb{R}$  is a scalar (so that  $e^{\Delta_t A_t}$  is a scalar and commutes with rotations), and  $\mathbf{R}_t$  is block-diagonal orthogonal/unitary, hence  $\mathbf{R}_t^{-1} = \mathbf{R}_t^\top$ .

Unrolling the recurrence with the convention that an empty product is the identity,

$$\mathbf{h}_t = \sum_{i=0}^t \left( \prod_{s=i+1}^t e^{\Delta_s A_s} \mathbf{R}_s \right) \mathbf{B}_i x_i. \quad (12)$$

Thus

$$\mathbf{y}_t = \mathbf{C}_t^\top \mathbf{h}_t = \sum_{i=0}^t \mathbf{C}_t^\top \left( \prod_{s=i+1}^t e^{\Delta_s A_s} \mathbf{R}_s \right) \mathbf{B}_i x_i. \quad (13)$$

Using unitarity property,

$$\prod_{s=i+1}^t \mathbf{R}_s = \left( \prod_{s=0}^t \mathbf{R}_s \right) \left( \prod_{s=0}^i \mathbf{R}_s \right)^{-1} = \left( \prod_{s=0}^t \mathbf{R}_s \right) \left( \prod_{s=0}^i \mathbf{R}_s^\top \right).$$

Since  $e^{\Delta_s A_s}$  are scalars, they commute with rotations; hence

$$\mathbf{y}_t = \sum_{i=0}^t \mathbf{C}_t^\top \left( \prod_{s=0}^t \mathbf{R}_s \right) \left( \prod_{s=i+1}^t e^{\Delta_s A_s} \right) \left( \prod_{s=0}^i \mathbf{R}_s^\top \right) \mathbf{B}_i x_i \quad (14)$$

$$= \left( \left( \prod_{s=0}^t \mathbf{R}_s^\top \right) \mathbf{C}_t \right)^\top \sum_{i=0}^t \left( \prod_{s=i+1}^t e^{\Delta_s A_s} \right) \left( \prod_{s=0}^i \mathbf{R}_s^\top \right) \mathbf{B}_i x_i. \quad (15)$$

Define the rotated parameters  $\bar{\mathbf{C}}_t := \left( \prod_{s=0}^t \mathbf{R}_s^\top \right) \mathbf{C}_t$  and  $\bar{\mathbf{B}}_i := \left( \prod_{s=0}^i \mathbf{R}_s^\top \right) \mathbf{B}_i$ . Then

$$\mathbf{y}_t = \bar{\mathbf{C}}_t^\top \sum_{i=0}^t \left( \prod_{s=i+1}^t e^{\Delta_s A_s} \right) \bar{\mathbf{B}}_i x_i. \quad (16)$$

Equivalently, introducing the rotated state  $\tilde{\mathbf{h}}_t := \left( \prod_{s=0}^t \mathbf{R}_s^\top \right) \mathbf{h}_t$ ,

$$\tilde{\mathbf{h}}_t = e^{\Delta_t A_t} \tilde{\mathbf{h}}_{t-1} + \bar{\mathbf{B}}_t x_t, \quad \mathbf{y}_t = \bar{\mathbf{C}}_t^\top \tilde{\mathbf{h}}_t, \quad (17)$$

$\square$



## C.3 PROOF OF PROPOSITION 4

**Proposition 4** (Rotary Embedding Equivalence with Trapezoidal Discretization). *Discretizing a complex SSM with the trapezoidal rule (Proposition 1) yields the recurrence*

$$\begin{aligned} \mathbf{h}_t &= \alpha_t \mathbf{h}_{t-1} + \beta_t \left( \prod_{i=0}^{t-1} \mathbf{R}_i^\top \right) \mathbf{B}_{t-1} x_{t-1} + \gamma_t \left( \prod_{i=0}^t \mathbf{R}_i^\top \right) \mathbf{B}_t x_t, \\ \mathbf{y}_t &= \left( \left( \prod_{i=0}^t \mathbf{R}_i^\top \right) \mathbf{C}_t \right)^\top \mathbf{h}_t. \end{aligned} \quad (9)$$

Here  $\mathbf{R}_t$  is the block-diagonal rotation matrix defined in Proposition 3.

## C.4 PROOF OF PROPOSITION 4

**Proposition 4** (Rotary Embedding Equivalence with Trapezoidal Discretization). *Discretizing a complex SSM with the trapezoidal rule (Proposition 1) yields the recurrence*

$$\begin{aligned} \mathbf{h}_t &= \alpha_t \mathbf{h}_{t-1} + \beta_t \left( \prod_{i=0}^{t-1} \mathbf{R}_i^\top \right) \mathbf{B}_{t-1} x_{t-1} + \gamma_t \left( \prod_{i=0}^t \mathbf{R}_i^\top \right) \mathbf{B}_t x_t, \\ \mathbf{y}_t &= \left( \left( \prod_{i=0}^t \mathbf{R}_i^\top \right) \mathbf{C}_t \right)^\top \mathbf{h}_t. \end{aligned} \quad (9)$$

Here  $\mathbf{R}_t$  is the block-diagonal rotation matrix defined in Proposition 3.

## C.5 PROOF OF PROPOSITION 4

**Proposition 4** (Rotary Embedding Equivalence with Trapezoidal Discretization). *Discretizing a complex SSM with the trapezoidal rule (Proposition 1) yields the recurrence*

$$\begin{aligned} \mathbf{h}_t &= \alpha_t \mathbf{h}_{t-1} + \beta_t \left( \prod_{i=0}^{t-1} \mathbf{R}_i^\top \right) \mathbf{B}_{t-1} x_{t-1} + \gamma_t \left( \prod_{i=0}^t \mathbf{R}_i^\top \right) \mathbf{B}_t x_t, \\ \mathbf{y}_t &= \left( \left( \prod_{i=0}^t \mathbf{R}_i^\top \right) \mathbf{C}_t \right)^\top \mathbf{h}_t. \end{aligned} \quad (9)$$

Here  $\mathbf{R}_t$  is the block-diagonal rotation matrix defined in Proposition 3.

*Proof.* We begin from the complex SSM (as in Prop. 2)

$$\begin{aligned} \dot{\mathbf{h}}(t) &= \text{Diag}(A(t) + i\theta(t)) \mathbf{h}(t) + (\mathbf{B}(t) + i\hat{\mathbf{B}}(t))x(t), \\ y(t) &= \text{Re}((\mathbf{C}(t) + i\hat{\mathbf{C}}(t))^\top \mathbf{h}(t)), \end{aligned}$$

where  $A(t) \in \mathbb{R}$  is a scalar and  $\theta(t), \mathbf{B}(t), \hat{\mathbf{B}}(t), \mathbf{C}(t), \hat{\mathbf{C}}(t) \in \mathbb{R}^{N/2}$ .

Recall from Prop. 5,

$$\mathbf{h}_t \approx e^{\Delta_t(A_t + i\theta_t)} \mathbf{h}_{t-1} + \int_{\tau_{t-1}}^{\tau_t} e^{(\tau_t - \tau)(A_t + i\theta_t)} (\mathbf{B}(\tau) + i\hat{\mathbf{B}}(\tau)) x(\tau) d\tau.$$

Applying Prop. 1 to the above integral, we get

$$\mathbf{h}_t = e^{\Delta_t(A_t + i\theta_t)} \mathbf{h}_{t-1} + \beta_t e^{i\Delta_t\theta_t} (\mathbf{B}_{t-1} + i\hat{\mathbf{B}}_{t-1}) x_{t-1} + \gamma_t (\mathbf{B}_t + i\hat{\mathbf{B}}_t) x_t, \quad (18)$$

wherem

$$\alpha_t := e^{\Delta_t A_t}, \quad \beta_t := (1 - \lambda_t) \Delta_t e^{\Delta_t A_t}, \quad \gamma_t := \lambda_t \Delta_t,$$

Since  $e^{\Delta_t(A_t + i\theta_t)} = \alpha_t e^{i\Delta_t\theta_t}$  and as shown in Prop. 2, multiplication by  $e^{i\Delta_t\theta_t}$  is a block-diagonal rotation in real coordinates, we get the real  $N$ -dimensional recurrence

$$\begin{aligned} \mathbf{h}_t &= \alpha_t \mathbf{R}_t \mathbf{h}_{t-1} + \beta_t \mathbf{R}_t \mathbf{B}_{t-1} x_{t-1} + \gamma_t \mathbf{B}_t x_t, \\ \mathbf{y}_t &= \mathbf{C}_t^\top \mathbf{h}_t, \end{aligned} \quad (19)$$

where  $\mathbf{R}_t = \text{Block}(\{R(\Delta_t \theta_t[i])\}_{i=1}^{N/2})$  where  $R(\Theta) = \begin{bmatrix} \cos \Theta & -\sin \Theta \\ \sin \Theta & \cos \Theta \end{bmatrix}$ , and projections

$\mathbf{B}_t = \begin{bmatrix} \mathbf{B}_t \\ \hat{\mathbf{B}}_t \end{bmatrix}$ ,  $\mathbf{C}_t = \begin{bmatrix} \mathbf{C}_t \\ -\hat{\mathbf{C}}_t \end{bmatrix}$ . Note that  $\mathbf{R}_t$  is orthogonal, so  $\mathbf{R}_t^{-1} = \mathbf{R}_t^\top$ .

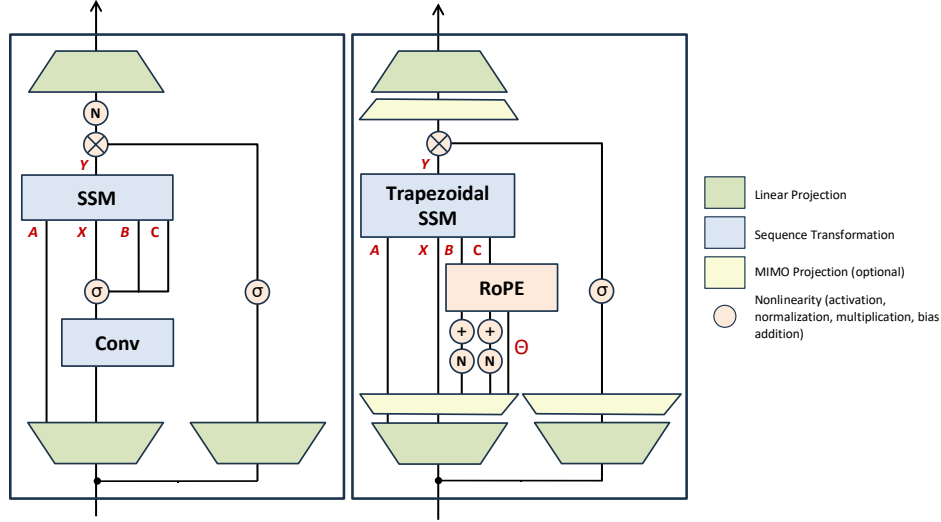


Figure 4: Contrasting Mamba-2 and Mamba-3 Architectures: Key updates include trapezoidal discretization, data-dependent RoPE embeddings, MIMO projections, QK normalization, and learnable biases.

We define the following,

$$\tilde{h}_t := \left( \prod_{s=0}^t \mathbf{R}_s^\top \right) h_t, \quad \bar{\mathbf{B}}_t := \left( \prod_{s=0}^t \mathbf{R}_s^\top \right) \mathbf{B}_t, \quad \bar{\mathbf{C}}_t := \left( \prod_{s=0}^t \mathbf{R}_s^\top \right) \mathbf{C}_t.$$

Left-multiplying equation 19 by  $\prod_{s=0}^t \mathbf{R}_s^\top$  and using  $\mathbf{R}_t^\top \mathbf{R}_t = \mathbf{I}$ ,

$$\begin{aligned} \tilde{h}_t &= \alpha_t \tilde{h}_{t-1} + \beta_t \bar{\mathbf{B}}_{t-1} x_{t-1} + \gamma_t \bar{\mathbf{B}}_t x_t, \\ y_t &= \bar{\mathbf{C}}_t^\top \tilde{h}_t. \end{aligned}$$

This is a vanilla scalar-transition SSM with data-dependent rotary embeddings absorbed into  $\mathbf{B}, \mathbf{C}$  via cumulative products of  $\mathbf{R}_s^\top$ .  $\square$

## D MIMO FOR MAMBA-3

With hindsight from Mamba and with inference in mind, we propose the following MIMO formulation:

**Mamba with MIMO** With a given batch, head, and sequence position  $t$ , consider the input  $\mathbf{U}_t \in \mathbb{R}^D$ . Also denote  $P, R \in \mathbb{N}$  as the head dimension and MIMO rank, respectively. We first obtain SSM parameters via a set of projections defined in terms of tensor contraction notation as follows:

$$\begin{aligned} \mathbf{B}_t &= \text{contract}(DNR, D \rightarrow NR)(\mathbf{W}_B, \mathbf{U}_t) & \mathbf{C}_t &= \text{contract}(DNR, D \rightarrow NR)(\mathbf{W}_C, \mathbf{U}_t), \\ \mathbf{X}'_t &= \text{contract}(PD, D \rightarrow P)(\mathbf{W}_{X'}, \mathbf{U}_t) & \mathbf{X}_t &= \text{contract}(PR, P \rightarrow PR)(\mathbf{W}_X, \mathbf{X}'_t), \end{aligned}$$

where  $\mathbf{W}_B, \mathbf{W}_C, \mathbf{W}_{X'}, \mathbf{W}_X$  are model parameters. Additionally, we obtain the residual term  $\mathbf{Z}_t$  in the same manner as  $\mathbf{X}_t$  with weights  $\mathbf{W}_{Z'}$  and  $\mathbf{W}_Z$ . The state update and the SSM output is then computed via the following MIMO SSM:

$$\mathbf{H}_t = a_t \mathbf{H}_{t-1} + \mathbf{B}_t \mathbf{X}_t^\top \in \mathbb{R}^{N \times P}, \quad \mathbf{Y}_t = \mathbf{H}_t^\top \mathbf{C}_t \in \mathbb{R}^{P \times R}.$$

The intermediate output  $\mathbf{Y}'_t$  is obtained via some residual function  $\phi$ ,  $\mathbf{Y}'_t \leftarrow \phi(\mathbf{Y}_t, \mathbf{Z}_t)$ . Finally, the layer output  $\mathbf{O}_t \in \mathbb{R}^D$  is computed via the following down projections:

$$\mathbf{O}'_t = \text{contract}(PR, R \rightarrow P)(\mathbf{W}_{O'}, \mathbf{Y}'_t) \quad \mathbf{O}_t = \text{contract}(P, PD \rightarrow D)(\mathbf{W}_O, \mathbf{O}'_t).$$

This formulation enhances the existing Mamba3 architecture by providing a lightweight parameterization that transforms the set of independent SISO SSMs within each head into a set of MIMO SSMs. Here, we note that the hardware-efficient chunking technique employed by Mamba2 for pretraining can be applied with little change, as the MIMO dimension  $r$  is orthogonal to the sequence dimension.

## E EXPERIMENTAL DETAILS

**Language Modeling** Our pretraining procedures follow that of Dao & Gu (2024)’s section D.2. All models at each scale follow the same procedure and were trained with bfloat16. The Mamba family of models were trained using the standard expand factor of 2 and a dstate of 128 and head dimension of 64. The Transformer baselines follows Dao & Gu (2024), and the Gated DeltaNet baselines follow (Yang et al., 2025a). We utilize the Llama-3.1 tokenizer (Grattafiori et al., 2024) for all models.

We utilize LM Evaluation Harness (Gao et al., 2024) to test the zero-shot language modeling capabilities of our pretrained model on LAMBADA (OpenAI version) (Paperno et al., 2016), HellaSwag (Zellers et al., 2019), PIQA (Bisk et al., 2019), Arc-Easy/Arc-Challenge (Clark et al., 2018), WinoGrande (Sakaguchi et al., 2019), and OpenBookQA (Mihaylov et al., 2018).

**Real-World and Synthetic Retrieval** For our real-world retrieval tasks, we evaluate on the common suite consisting of SWDE (Arora et al., 2025b), SQUAD (Rajpurkar et al., 2018), FDA (Arora et al., 2025b), TriviaQA (Joshi et al., 2017), NQ (Kwiatkowski et al., 2019), and DROP (Dua et al., 2019). We utilize the cloze-formatted version of the aforementioned tasks provided by Arora et al. (2025b; 2024), as the original datasets are in a question-answering format, making it challenge for solely pretrained models. All tasks were truncated to match the training context length. The synthetic NIAH tasks (Hsieh et al., 2024) were also run with LM Evaluation Harness.

**State-Tracking Synthetics** Training follows a sequence length curriculum that progresses from 3 -40 to 160, evaluated at 256. Each curriculum runs for  $10^4$  steps with batch size 256. We use 1 layer models for Parity and 3 layer models for Modular-arithmetic tasks. The state size is chosen to be 64, and we sweep  $d_{\text{model}} \in \{32, 64\}$  and 8 learning rates logarithmically spaced between  $10^{-4}$  and  $10^{-2}$ , reporting the best validation accuracy.

## F ADDITIONAL EXPERIMENTAL RESULTS

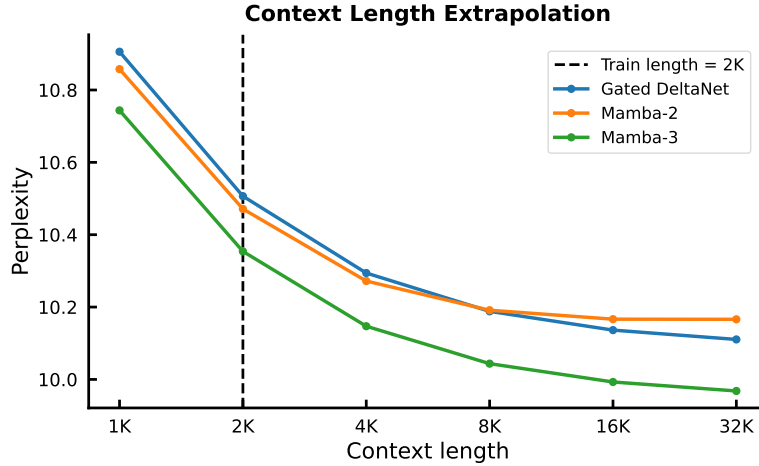


Figure 5: Pretrained 1.5B models’ performance on the held-out FineWeb-Edu test set at varying context lengths. Mamba-3 exhibits strong length extrapolation while Mamba-2 falters at longer contexts.

Table 5: Downstream language modeling evaluations on parameter-matched pretrained models, including Mamba-3 MIMO. Mamba-3 MIMO’s average accuracy on all tasks is more than 1 percentage point better than the next best (Mamba-3 SISO).

Model	FW-Edu ppl ↓	LAMB. ppl ↓	LAMB. acc ↑	HellaS. acc.n ↑	PIQA acc ↑	Arc-E acc ↑	Arc-C acc.n ↑	WinoGr. acc ↑	OBQA acc.n ↑	Average acc ↑
Transformer-440M	13.03	21.2	41.7	50.5	69.9	67.6	34.6	<b>56.7</b>	26.0	49.6
Gated DeltaNet-440M	13.12	19.0	40.4	50.5	70.5	67.5	34.0	55.3	25.8	49.1
Mamba-2-440M	13.00	19.6	40.8	51.7	70.6	68.8	35.0	54.1	26.0	49.6
<b>Mamba-3-440M</b>	<u>12.87</u>	19.6	40.2	<u>51.7</u>	<b>71.9</b>	<u>68.9</u>	34.4	55.8	26.0	<u>49.8</u>
<b>Mamba-3-MIMO-440M</b>	<b>12.72</b>	17.1	43.4	<b>52.8</b>	<u>70.8</u>	<b>69.6</b>	<b>35.6</b>	<u>56.3</u>	<b>28.4</b>	<b>51.0</b>

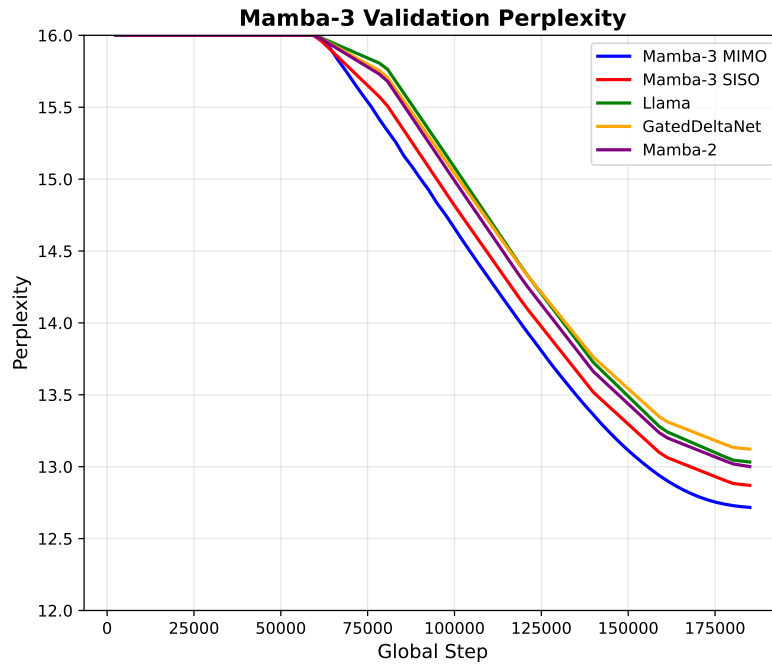


Figure 6: Mamba-3 demonstrates superior performance compared to strong baselines like Mamba-2, Llama, and Gated Deltanet. These are 440M models, evaluated on FineWeb-Edu and 100B tokens.