

Transformers to SSMs: Distilling Quadratic Knowledge to Subquadratic Models

Aviv Bick*¹³, Kevin Y. Li*², Eric P. Xing²⁴, J. Zico Kolter², and Albert Gu²³

¹Computer Science Department, Carnegie Mellon University

²Machine Learning Department, Carnegie Mellon University

³Cartesia.ai

⁴MBZUAI

abick@cs.cmu.edu, kyl2@cs.cmu.edu

Abstract

Transformer architectures have become a dominant paradigm for domains like language modeling but suffer in many inference settings due to their quadratic-time self-attention. Recently proposed subquadratic architectures, such as Mamba, have shown promise, but have been pretrained with substantially less computational resources than the strongest Transformer models. In this work, we present a method that is able to distill a pretrained Transformer architecture into alternative architectures such as state space models (SSMs). The key idea to our approach is that we can view both Transformers and SSMs as applying different forms of mixing matrices over the token sequences. We can thus progressively distill the Transformer architecture by matching different degrees of granularity in the SSM: first matching the mixing matrices themselves, then the hidden units at each block, and finally the end-to-end predictions. Our method, called MOHAWK, is able to distill a Mamba-2 variant based on the Phi-1.5 architecture (Phi-Mamba) using only 3B tokens and a hybrid version (Hybrid Phi-Mamba) using 5B tokens. Despite using less than 1% of the training data typically used to train models from scratch, Phi-Mamba boasts substantially stronger performance compared to all past open-source non-Transformer models. MOHAWK allows models like SSMs to leverage computational resources invested in training Transformer-based architectures, highlighting a new avenue for building such models.

1 Introduction

Large language models based upon Transformer architectures have become a staple of natural language processing but suffer from their reliance on quadratic self-attention – the need to compute inner products between tokens at all positions up to the context length. This has motivated the development of several alternative subquadratic models, either approximations of self-attention (Katharopoulos et al. 2020) or entirely different architectures, such as state space models (SSMs) (A. Gu and Dao 2023; A. Gu, Goel, and Ré 2022; Peng et al. 2023; Sun et al. 2023). Training strong subquadratic models such as SSMs can benefit the community through their cheaper finetuning and inference costs; however, they have not benefited from the same amount of community effort in the form of training and compute as for Transformers. This raises a natural question: is it possible to leverage the vast amounts of resources that have been invested in training quadratic-time Transformers and use these models to produce stronger alternative models, such as state-space models?

In this paper, we present an approach for training subquadratic state-space models (specifically from the class of Mamba SSMs (A. Gu and Dao 2023)) through the distillation of different elements of a pretrained Transformer model. The key intuition is viewing both Attention and SSMs as sequence transformations that mix different token embeddings by applying different classes of matrices across them. Sequence model *architectures* can then be factored into separate (i) sequence mixing and (ii) channel mixing blocks, e.g., a Transformer is composed of Attention (sequence mixer) and MLP (channel mixer) blocks. Using this breakdown, we can separately distill the *mixing* elements of each model explicitly at different

* Authors contributed equally to this work.

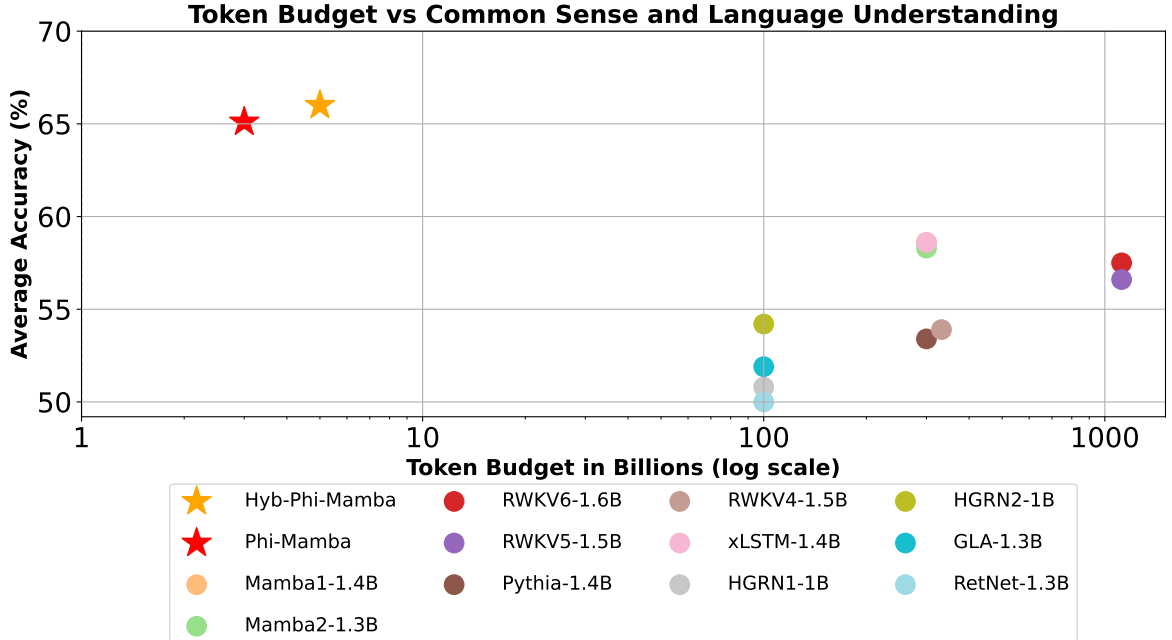


Figure 1: Plot of trained token budget to averaged accuracy on Winogrande, Arc-E, Arc-C, PIQA, and Hellaswag on various open-source models (mainly non-Transformer-based models). Our model (Phi-Mamba) uses more than 33× less token budget to achieve 5% higher average accuracy than the next best model.

levels of granularity. Specifically, we propose a three-phase distillation process that progressively targets higher levels of supervision from the teacher model: (1) a matrix orientation phase that aligns the sequence transformation matrices themselves; (2) a hidden-state distillation that aligns the hidden-state representations of each individual layer of the network without sacrificing preexisting learned representations; and (3) an end-to-end training phase with weight transfer that finally distills the final output of the network using only a fraction of training data. We term our approach **MOHAWK** after these three stages (**M**atrix **O**rientation, **H**idden-State **A**lignment, **W**eight-Transfer and **K**nowledge **D**istillation).

We apply our approach to a modified instantiation of the Mamba-2 architecture (Dao and A. Gu 2024), termed Phi-Mamba, which is aimed at more directly corresponding to the different architectural blocks of the Phi-1.5 language model (Gunasekar et al. 2023) – a very strong Transformer model at the 1.3B parameter scale. Using our approach, the Phi-Mamba model achieves performance on benchmarks *stronger than any previous Mamba model of similar size*. Although performance still lags behind that of the base Phi-1.5 model on these benchmarks, the model is distilled with only 3.0B tokens, less than 1% of the data used to train either the previously best-performing Mamba models and 2% for the Phi-1.5 model itself. For instance, our Phi-Mamba achieves a 71.7% accuracy on the Winogrande dataset, compared to the pretrained Mamba-2 model’s 60.9% accuracy, and 44.1% accuracy on the ARC-C dataset, compared to the Mamba-2’s 33.3% accuracy.

Moreover, by retaining only four attention layers and substituting the remaining 20 layers with Mamba, our hybrid model attains an average performance of 66.0% on select downstream evaluations, compared to Phi-1.5’s 67.2%. Interestingly, our hybrid Phi-Mamba surpasses or closely matches Samba (Ren et al. 2024) on multiple benchmark tasks, even though Samba was trained using Phi-2’s dataset, contains more parameters, and features 3× the number of attention layers. Given that we utilized the lower quality C4 dataset (Raffel et al. 2023) for distillation, it indicates that distilling from the Transformer model (Phi-1.5 from Gunasekar et al. (2023)) *even without its original training data* can outperform training from scratch *on the same/comparable training data*.

Our results highlight the benefit of our three-phase distillation approach: we show in ablation experiments that each phase is highly beneficial for the eventual final performance of the model and that, e.g., *only* attempting to directly distill the Phi-1.5 model (i.e., Phase 3 alone) substantially underperforms the full MOHAWK method. Moreover, our findings emphasize the benefits of state-space models while training on fewer than 100× tokens than the original pretrained Mamba model.

2 Related Work

Sequence Models. State-of-the-art autoregressive language models have been pretrained on massive amounts of data, resulting in models that exhibit extensive downstream capabilities, such as zero-shot translation and long-range reasoning (Brown et al. 2020; Gunasekar et al. 2023; Touvron et al. 2023). Recent work has focused on addressing the quadratic complexity of Transformers by developing subquadratic alternatives based on RNN (Beck et al. 2024; Peng et al. 2023), SSM (A. Gu and Dao 2023; Sun et al. 2023), and linear attention mechanisms (Dao, Fu, et al. 2022; Katharopoulos et al. 2020; Liu, Zaharia, and Abbeel 2023; Qin, S. Yang, et al. 2024; S. Yang et al. 2024), highlighting the importance of efficient sequence models in the era of large-scale autoregressive language models.

In addition, to combine different capabilities while maintaining efficiency, hybrid models that integrate attention mechanisms with subquadratic methods have been proposed (Fu et al. 2023; Lieber et al. 2024; Ren et al. 2024; Z. Wang et al. 2024). These models typically feature a limited number of attention layers, thus maintaining the quadratic complexity at a relatively low factor.

SSM Architectures. GSS was the pioneer in integrating SSMs into gated neural network architecture for language modeling (Mehta et al. 2022). H3, inspired by the combination of S4 and linear attention (Katharopoulos et al. 2020), employs SSMs with shift and diagonal matrices and multiplicative operations on input projections, extending this formulation to broader recurrences, a foundation for subsequent architectures (Fu et al. 2023). Selective S4 incorporates S4 as a black box that generates a binary mask applied to the input, an architectural modification akin to gating mechanisms (J. Wang et al. 2023). Mamba (A. Gu and Dao 2023), combines the H3 block with the ubiquitous MLP block of modern neural networks by interleaving them, resulting in a more powerful architecture. The Mamba-2 block simplifies the Mamba block by removing sequential linear projections; the SSM parameters A, B, C are produced at the beginning of the block rather than as a function of the input X of the SSM. Finally, Mamba-2 (Dao and A. Gu 2024) was introduced as a variant of Mamba which leverages the structured state space duality (SSD). The Mamba-2 core layer is 2-8x faster than Mamba’s selective SSM while continuing to outperform Transformers in language modeling.

Distillation. Knowledge distillation can be used to transfer knowledge from a large teacher model to a smaller student model, resulting in a more efficient model that retains the performance of the teacher model (Hinton, Vinyals, and Dean 2015). Distillation has been applied to various language modeling tasks, such as text generation (Chen et al. 2020; Haidar and Rezagholizadeh 2019), machine translation (Hahn and H. Choi 2019; Tan et al. 2019; Zhou, Neubig, and J. Gu 2021), and question-answering system (Hu et al. 2018; Z. Yang et al. 2019).

Distillation in language models has been largely focused on *compression*: turning a larger pretrained Transformer into a smaller one by utilizing the weights of the teacher model (Jha et al. 2023; W. Wang et al. 2020; Xia et al. 2023). Some of the techniques proposed look similar to ours; for example, W. Wang et al. (2020) match attention matrices in a step similar to our matrix orientation, and Liang et al. (2023) align outputs of each block (i.e., the hidden states). However, these differ in subtle and important ways because of our setting; for example, the former uses a different loss function than us that relies on softmax attention, and the latter is an end-to-end objective while our hidden state alignment occurs completely independently block-per-block. Consequently, prior work has observed that combining these objectives does not actually help and even might hurt distillation (Jha et al. 2023), whereas we show that our techniques all significantly help improve the student model.

A smaller body of work has focused on our objective of distilling *across architectures*, in particular, turning a pretrained Transformer into a different architecture (usually a recurrent model of some form) of the same size. Kasai et al. (2021) first tried turning a pretrained softmax attention into a linear attention by directly transferring weights and continuing fine-tuning; a similar approach was taken by concurrent work (Mercat et al. 2024). Recently, Zhang et al. (2024) also proposed distilling into linear attention by first matching attention matrices. Our approach differs by using a more constrained loss function that works beyond linear attention; incorporating more fine-grained alignment (e.g., the hidden state alignment step); and using recent, more expressive classes of efficient student models (Mamba-2), which we show are significantly easier to distill (Table 7).

3 Background and Overview

To facilitate a clear understanding of our distillation approach, we start with the necessary background and definitions. An overview of the Mamba-2 architecture, which forms the foundation of our Phi-Mamba model, is also provided.

3.1 Matrix Mixers

Following Dao and A. Gu (2024), we refer to an equivalent function that represents the input and output of a sequence model as a *sequence transformation* or a *sequence mixer*. Formally,

Definition 1 (Sequence Transformation). *We use the term sequence transformation to refer to a parameterized map on sequences $Y = f_{\theta}(X)$ where $X, Y \in \mathbb{R}^{(T,P)}$ and θ is an arbitrary collection of parameters. T represents the sequence or time axis; subscripts index into the first dimension, e.g. $X_t, Y_t \in \mathbb{R}^P$.*

To put it differently, sequence mixers combine tokens at various time steps, facilitating the model’s comprehension of temporal information and interactions. Sequence transformations form the foundation of deep sequence models, being integral components of neural network frameworks such as Transformers. A particular family of sequence transformations can be represented by $Y = MX$ for a matrix $M \in \mathbb{R}^{(T,T)}$, which we refer to as a *sequence transformation matrix* or *matrix mixer*.

An example of such a matrix mixer is the vanilla self-attention, $\text{Softmax}(QK^T)$, which is applied to the input-dependent V resulting in the familiar $\text{Softmax}(QK^T)V$. Similarly, Linear Attention (Katharopoulos et al. 2020) has a sequence transformation matrix of the form K^T . In addition, we can easily obtain their causal variants by multiplying by L , a lower triangular matrix filled with 1s, to obtain $L \circ \text{Softmax}(QK^T)$ and $L \circ QK^T$, respectively. Another example is a Toeplitz matrix T used to perform discrete convolution on input X , resulting in TX (Qin, Han, et al. 2023).

A naive approach to computing the output of a sequence transformation is to multiply the input sequence X by the matrix M . However, this approach has a time complexity of $O(T^2)$, which is prohibitive for long sequences. Subquadratic sequence transformations, such as Mamba-2, have been developed to address such inefficiencies through structured matrix multiplication.

3.2 Mamba-2

Mamba-2 (Dao and A. Gu 2024), a type of structured state space models (SSMs) (A. Gu 2023; A. Gu, Goel, and Ré 2022), was recently introduced. Similarly to the original Mamba model (A. Gu and Dao 2023), Mamba-2 uses a time-varying state-space model which can selectively focus on or ignore inputs due to its input-dependent parameterization of the system components. The time-varying SSM is defined as follows:

$$\begin{aligned} h_{t+1} &= A_t h_t + B_t x_t \\ y_t &= C_t h_t \end{aligned} \tag{1}$$

Here, B_t and C_t are input-dependent projections of the system, as in Mamba-1; however, A_t is the identity matrix I multiplied by a scalar α_t . The above formulation also differs from the previous one by treating the underlying sequence as originating from a discrete signal instead of a continuous one and therefore omits the sampling component Δt from the original Mamba model.

Importantly, Mamba-2 draws a new connection between SSMs and Transformers, termed *Structured State Space Duality (SSD)*, which shows that a special case of SSMs can be viewed as a form of causal linear attention. In particular, fixing $A_t = I$ (a further restriction of Mamba-2 to $\alpha_t = 1$) results in the formulation of causal linear attention (Katharopoulos et al. 2020) with the matrices B and C representing the projections of the key and the query, respectively, while the input projection X corresponds to the projection of the value.

Mamba-2 as a matrix sequence transformation. Inspired by the aforementioned connection between SSMs and Transformers, Dao and A. Gu (2024) shows that Mamba-2’s SSD mixer family is equivalent to sequentially-semi-separable

matrices (Chandrasekaran et al. 2002). Formally, the SSD mixer family can be represented as:

$$\begin{cases} h_{t+1} = \alpha_t \cdot I h_t + \mathbf{B} x_t \\ y_t = \mathbf{C} \cdot h_t \end{cases} \Rightarrow \begin{bmatrix} \alpha_1 & 0 & 0 & \cdots & 0 \\ \alpha_{2:1} & \alpha_2 & 0 & \cdots & 0 \\ \alpha_{3:1} & \alpha_{3:2} & \alpha_3 & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \alpha_{n:1} & \alpha_{n:2} & \alpha_{n:3} & \cdots & \alpha_n \end{bmatrix} \circ (\mathbf{C} \cdot \mathbf{B}^\top) \cdot X \quad (2)$$

where $\alpha_{t,i} = \alpha_{t-1} \cdot \alpha_{t-2} \cdots \alpha_i$. An interesting observation is that the Mamba-2 architecture can be viewed as a causal linear attention with a *learnable causal mask*.

The Mamba-2 block. To enhance the effectiveness of the above Mamba-2 matrix mixer (Equation (2)), Dao and A. Gu (2024) design the Mamba-2 block, a modified version of the Mamba-1 block (A. Gu and Dao 2023). They added parallel parameter projections, where $\mathbf{A}, \mathbf{X}, \mathbf{B}, \mathbf{C}$ are produced in parallel, reducing parameters and supporting tensor parallelism, and an extra normalization layer before the final output projection to address instabilities in training larger models.

Although we have made additional modifications to the Mamba-2 block, they remain quite similar. Therefore, for a visual representation of the Mamba-2 block, refer to Figure 2. Note that we have reverted the introduced normalization layer before the final output projection and have also discarded the nonlinear activation after the convolution operation found in both the Mamba-1 and Mamba-2 blocks.

4 Methods

Throughout this section, we will describe each phase of MOHAWK. Specifically, we will cover the stages of matrix orientation, hidden-state alignment, and knowledge distillation, all three of which are crucial for developing an effective student model from the pretrained Transformer model. Unlike traditional distillation techniques, the student model retains the overall architecture of the teacher model, differing only in the replacement of the attention matrix mixer with a subquadratic alternative. We will progressively unveil our architecture, Phi-Mamba, along with the specifics of its distillation process. This section concludes with an in-depth description of the Phi-Mamba architecture and its hybrid version, which surpasses the performance of other subquadratic matrix mixers. Further examinations of the effectiveness of the method and ablation studies are discussed in Section 5.

For clarity, the term *block* refers to the repeating components that form the end-to-end model. The blocks are composed of *layers*, such as the self-attention layer (including projections), the SSM layer (including the mixer and convolution), and the convolutional layer. In this manner, many Transformer models, like Llama (Touvron et al. 2023), are viewed as a stack of alternating self-attention and MLP blocks, whereas the Phi and Phi-Mamba models are comprised of Phi blocks that have parallel Attention/SSM and MLP blocks.

4.1 Stage 1: Matrix Orientation

The first stage of MOHAWK aims to align the student matrix mixer with the teacher’s self-attention matrix. Achieving this alignment is a two-step process: first, at every mixing layer, the student components preceding the matrix mixer are set to match the teacher’s components. This ensures that each layer’s input undergoes the same transformation up to the matrix mixer section. Consequently, the only variation from the input to the mixing process is the matrix calculation. We then minimize the distance between the matrix mixer, e.g., the self-attention matrix and the materialized SSM matrix (2), of each layer within the student and teacher models:

$$\min_{\phi} \|\text{TeacherMixer}(\mathbf{u}) - \text{StudentMixer}_{\phi}(\mathbf{u})\|_F \quad (3)$$

where ϕ denotes the parameters within the student’s sequence mixing layer, and \mathbf{u} indicates any arbitrary input. In our experimental setup, \mathbf{u} was chosen as the output from the teacher model’s preceding layer to better mimic the input distribution to the layer. This stage ensures that the student and teacher models have roughly similar mixing layers and sets the foundation for the subsequent stages of the distillation process. In particular, this stage can be done in parallel across all the student layers, as the inputs to the student and teacher blocks are identical.

For Mamba-2, we begin by setting the convolution to an identity function, effectively nullifying its initial impact. This results in the computation of the semi-separable matrix being the sole distinction between the layers. We then proceed to minimize the distance between the two matrix mixers: the semiseparable scalar identity and the attention matrix (see Figure 2). Figure 3 demonstrates the importance of this stage in the distillation process. Furthermore, Table 6 shows that the Mamba-2 matrix mixer is more expressive than popular alternatives and can closely approximate the self-attention matrix of various data samples across all layers of a Transformer model through gradient descent, solidifying it as a strong sequence mixer.

4.2 Stage 2: Hidden-State Alignment

Following the optimization of Equation (3), we must still address the differences between the outputs of the student and teacher blocks. To achieve this, we further align the components of the two blocks using initialization and distillation. Specifically, our goal is to match each student and teacher mixing blocks by minimizing the L2 norm of their output (e.g., the entire Mamba block with the self-attention block):

$$\min_{\phi} \|\text{AttnBlock}(\mathbf{u}) - \text{StudentMixerBlock}_{\phi}(\mathbf{u})\|_2 \tag{4}$$

where similar to Section 4.1, ϕ represents student’s block parameters, and \mathbf{u} is an input. Once again, this stage can be done in parallel across all the student layers.

In the case of Mamba-2, we modify the remaining components to be identical to the Phi-1.5’s Attention block, so that the overall functionality is preserved from Stage 1. Concretely, we initialize the gate (see Figure 2) to a constant value of 1 to “open” the gate, canceling its initial effect. In addition, we remove the normalization prior to the output projection, as it cannot be set to align with the Attention block. We then minimize the distance between the output of the Mamba-2 block and the output of the teacher’s self-attention block. Our analysis indicates that the distance between the Mamba-2 block and the self-attention block is strongly correlated with the model’s ability to learn the teacher’s distribution, as shown in Table 6. Furthermore, Figure 3 shows that a better independent alignment of the student and teacher blocks results in performance improvements, highlighting the importance of this stage in the distillation process.

4.3 Stage 3: Weight-Transfer and Knowledge Distillation

The final stage of the distillation process aims to fine-tune the student model to match the performance of the teacher model. Although each student mixing block is aligned with its corresponding teacher mixing block, discrepancies are still present between consecutive blocks throughout the network. To bridge these gaps and address the remaining components of the language model, we transfer the remaining weights of the teacher model to the student’s respective components. For Phi-Mamba, this involves the token embedding, the final layer normalization, the Language Model head, and the MLP and input norm at each block (see Figure 2). We then fine-tune the complete end-to-end student model under teacher supervision. Concretely, we use a distillation loss to encourage the student model to mimic the distribution of the teacher model’s logits, also known as knowledge distillation (Hinton, Vinyals, and Dean 2015):

$$\min_{\phi} \mathcal{L}_{\text{CE}}(\text{TeacherModel}(\mathbf{x}), \text{StudentModel}_{\phi}(\mathbf{x})) \tag{5}$$

where \mathbf{x} is the input tokens to the models.

It has been hypothesized that much of the information stored in language models resides in MLP blocks (Niu et al. 2024). To utilize the work already done pretraining the teacher, MOHAWK adjusts the structure of the student blocks to utilize the MLP in the same way as the teacher model, effectively swapping the teacher’s matrix mixer with that of the student.

Interestingly, during this step, the MLP weights *can be kept frozen* while keeping the model performant. This showcases Mamba-2’s powerful expressiveness crucial for replacing Attention, cuts the number of trained parameters by more than half, and, in larger models, helps prevent the student model from experiencing catastrophic forgetting of the teacher model’s information. We validate Mamba-2’s ability to do so in Table 8.

4.4 Phi-Mamba architecture

Combining the three stages of MOHAWK, we introduce the *Phi-Mamba* architecture, which merges the Mamba-2 model of Dao and A. Gu (2024) with the Phi-1.5 Transformer model of Gunasekar et al. (2023). It consists of a stack of Phi-

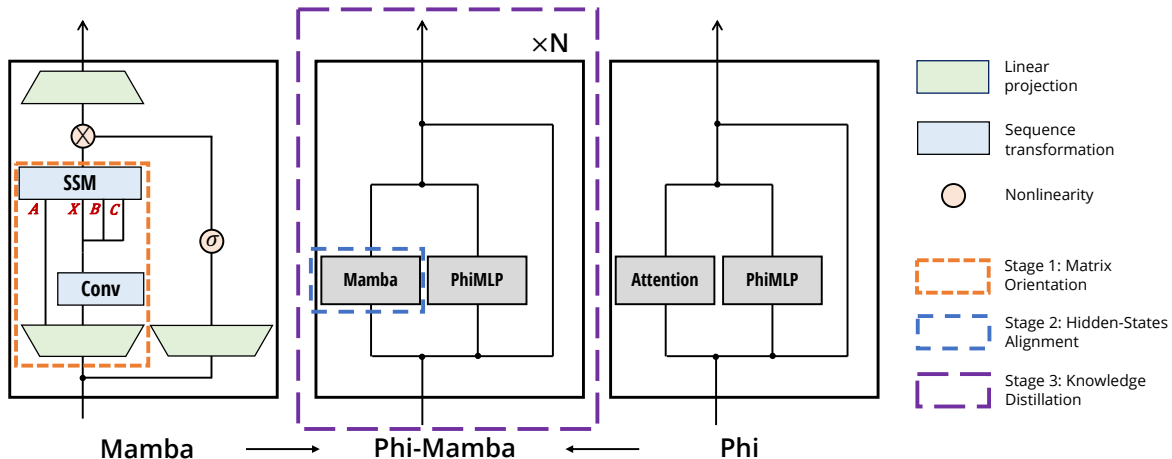


Figure 2: The Phi-Mamba architecture consists of a stack of blocks, each of which contains a Mamba block and an MLP block. The Mamba block is a simplified version of the Mamba-2 block (Dao and A. Gu 2024) that omits the non-linear activation function after the convolutional operation, as well as the layer normalization present before the output projection, so that the parts of the model outside the matrix mixer can be transferred from the teacher model. The MOHAWK distillation process involves progressively matching fine-to-coarse parts of the model to the corresponding part of the teacher model: (1) the mixer mixer itself (2) the full Mamba vs. Attention blocks, and (3) the end-to-end model.

Mamba blocks (Figure 2), initialized and distilled as described in previous sections. Additionally, we introduce the *Hybrid-Phi-Mamba* variant, which retains 4 layers of attention of Phi-1.5, effectively leveraging the strengths of both sequence mixers.

Overall, the Phi-Mamba architecture, as depicted in Figure 2, differs from the vanilla Mamba-2 architecture by modifying the structure of the SSM matrix mixer, removing components from the SSM block and incorporating dense layers from the teacher model. In particular, each Mamba-2 block was modified by removing post-convolution activation and pre-output projection normalization, while setting the gate and convolution to be identity functions. Interestingly, although these components were found to be beneficial for performance when Mamba-2 was trained from scratch (Dao and A. Gu 2024), we find that they are unnecessary for our distillation process.

Two key changes were made to the Mamba-2 matrix mixer. The first was converting the SSM head structure from multi-value to multi-head, much like the multi-head attention mechanism found in Transformers (Vaswani et al. 2023), enabling the independent distillation of each Transformer head into a Mamba head. Moreover, we handle the sequence mixer as entirely discrete-time by making the A matrix a projection of the input and eliminating the Δ discretization parameter. Although this formulation slightly differs from Mamba-2, the original algorithm can still be applied as a black-box method (refer to Appendix B).

5 Empirical Validation

We start by examining in Section 5.1 downstream evaluation scores of our MOHAWK-distilled Phi-Mamba-1.5B and Hybrid-Phi-Mamba-1.5B, empirically showing that they outperform all previous subquadratic and hybrid models, respectively, while having better time and memory complexities.

Next, Sections 5.2, 5.3, and 5.4 analyze our three-stage framework in reverse order of their introduction, disentangling the compounding effects of MOHAWK on the transfer of learned representations to the student model. Additionally, to form a baseline that mirrors the Phi-Mamba distillation process in ideal conditions, we employed MOHAWK to distill a Phi-1.5 into another Phi-1.5, transferring all weights except Attention layers, which were initialized from scratch. The specifications of our final Phi-Mamba model distilled using MOHAWK are provided in Section 5.5.

Section 5.6 outlines the architecture selected for the hybrid Phi-Mamba, discusses the ablations regarding the number and placement of interleaved attentions, and tackles a limitation potentially caused by the distillation process.

Table 1: Downstream evaluation results for full methods, comparing Phi-Mamba against open-source models of similar sizes pretrained on standard language modeling corpuses. Phi-Mamba attains performance close to the teacher model and better than all pretrained models, while using less than 1% of the training data.

| MODEL | TOKENS / DATASET | WINO.G. | ARC-E | ARC-C | PIQA | HELLAS. | LAMB. | AVG. \uparrow |
|-----------------------|----------------------|-------------|-------------|-------------|-------------|-------------|-------------|-----------------|
| Phi-1.5-1.3B | 150B / unknown | 73.4 | 75.6 | 48.0 | 76.6 | 62.6 | 53.4 | 64.9 |
| Phi-Mamba-1.5B | 3.0B / C4 | 71.7 | 74.0 | 44.1 | 75.5 | <u>60.2</u> | 50.1 | 62.6 |
| Mamba-1-1.4B | 315B / The Pile | <u>61.5</u> | <u>65.5</u> | 32.8 | <u>74.2</u> | 59.1 | 64.9 | <u>59.7</u> |
| Mamba-2-1.3B | 315B / The Pile | 60.9 | 64.3 | 33.3 | 73.2 | 59.9 | <u>65.7</u> | 59.6 |
| Finch-1.6B | 1.1T / RWKV World v2 | 59.4 | 64.2 | <u>34.1</u> | 72.6 | 57.3 | 66.8 | 59.1 |
| xLSTM-1.4B | 300B / SlimPajama | 60.6 | 64.3 | 32.6 | 74.6 | 60.9 | 57.8 | 58.5 |
| Eagle-1.5B | 1.1T / RWKV World v2 | 59.1 | 64.3 | 33.5 | 71.1 | 55.0 | <u>65.7</u> | 58.1 |
| Pythia-1.4B | 300B / The Pile | 57.3 | 60.6 | 26.0 | 71.1 | 52.1 | 61.6 | 54.8 |
| RWKV4-1.5B | 330B / The Pile | 54.6 | 60.5 | 29.4 | 72.4 | 52.5 | 56.4 | 54.3 |
| DeltaNet-1.3B | 100B / SlimPajama | 53.6 | 57.2 | 28.3 | 71.2 | 50.2 | 48.9 | 51.6 |
| GLA-1.3B | 100B / SlimPajama | 53.9 | 57.2 | 26.6 | 71.8 | 49.8 | 46.9 | 51.0 |

Lastly, Section 5.7 dives deeper into Mamba-2’s capability to learn interactions similar to that of Self-Attention. Section 5.7.1 examines the extent to which a Mamba-2 sequence transformation can approximate a self-attention matrix, and Section 5.7.2 examines whether this capability is evident in a comprehensive language model such as Phi-1.5.

5.1 Final Results

We empirically validate that our framework, MOHAWK, is able to achieve better performance on various downstream benchmarks compared to previous subquadratic models of similar size. We distill the Phi-1.5-1.3B model into Phi-Mamba-1.5B as well as Hybrid-Phi-Mamba-1.5B. Our final Phi-Mamba model is distilled on 3 billion tokens (distributed as 80M in Stage 1, 160M in Stage 2, and 2.76B tokens in Stage 3 as described in Section 5.5) from the C4 dataset, with a sequence length of 2048. This constitutes less than 1% of the resources used by many top-performing subquadratic open-source models (e.g., the open-source Mamba and Mamba-2 models, which are pretrained on 315 billion tokens). The Hybrid-Phi-Mamba-1.5B is distilled on a budget of 5 billion tokens from the same dataset.

Table 1 and Table 2 present a comprehensive breakdown of downstream evaluation results for our models and multiple baselines on a standard set of commonsense reasoning and language understanding tasks: WinoGrande (Sakaguchi et al. 2021), HellaSwag (Zellers et al. 2019), PIQA (Bisk et al. 2020), ARC-challenge and ARC-easy (Clark et al. 2018), and LAMBADA (Paperno et al. 2016). Figure 1 shows the performance versus the training cost of Phi-Mamba and Hybrid-Phi-Mamba compared to many open-source baselines from the literature at similar model sizes.

For the remainder of this section, we will analyze the impact of the 3 stages of MOHAWK one by one. Throughout the experiments detailed in this section, we use the AdamW optimizer with $\beta = (0.9, 0.95)$, a weight decay of 0.1, and a learning rate of 1×10^{-4} , combined with a Warmup-Stable-Decay (WSD) scheduler featuring 10% warmup and 10% decay. The training law figures and the final Phi-Mamba model use the regime detailed in Appendix A.

5.2 Stage 3 (Weight-Transfer and Knowledge Distillation)

As described in Section 4.3, this phase employs a simple end-to-end distillation of teacher-model logits. It leverages the alignment among all sequence mixers and successive blocks to jointly fine-tune all components of the network. Experiments shown in Table 3 highlight the relevance of implementing this end-to-end alignment, with all three architectures achieving their highest scores only after this phase. Predictably, the impact of end-to-end alignment varies by architecture: models with more mixing layers similar to the teacher model see a reduced importance of this phase.

Stage 3 is the only stage in MOHAWK that trains the student model end-to-end and can be seen as the “main” stage. Many distillation methods employ only this stage; however, Table 3 shows that using only end-to-end knowledge distillation is less than ideal. Although it is slightly advantageous to use only Stage 3 compared to only Stage 2 for both Phi-Mamba and Hybrid-Phi-Mamba, there is a significant gap between using only Stage 2 versus using Stage 2 + 3.

Table 2: Evaluation results show the performance of Hybrid-Phi-Mamba in downstream tasks, compared with similar-sized open-source models pretrained on standard language modeling data. Hybrid-Phi-Mamba utilizes under 3% of the training dataset and employs over 3x fewer attention layers (w represents the local window size of SWA). Both Samba and Mamba-SWA-MLP (Ren et al. 2024) stack layers of Mamba, Attention, and MLPs and are the only hybrid architectures of approximately 1.5B in size that we are aware of. An evaluation for Samba on the Lambda dataset was not available, hence it has been excluded.

| MODEL | # ATTNS | WINO.G. | ARC-E | ARC-C | PIQA | HELLAS. | AVG. \uparrow |
|--------------------------|-------------------|---------|-------------|-------------|------|-------------|-----------------|
| Phi-1.5-1.3B | 24 | 73.4 | 75.6 | 48.0 | 76.6 | 62.6 | 67.2 |
| H. Phi-Mamba-1.5B | 4 | 72.0 | 75.3 | 45.8 | 76.5 | 60.6 | 66.0 |
| Mamba-SWA-MLP-1.6B | 18 ($w = 2048$) | 73.7 | 76.6 | 46.1 | 76.5 | 49.7 | 64.5 |
| Samba-1.7B | 12 ($w = 2048$) | 72.9 | 79.2 | 48.2 | 77.1 | 49.7 | 65.4 |

Table 3: MOHAWK distillation was performed on the following models: (1) Phi-Mamba-1.5B, (2) Hybrid-Phi-Mamba-1.5B, and (3) Phi-1.5-1.3B. The teacher model for all three architectures was Phi-1.5. “Stages Applied” details which of the three MOHAWK stages was carried out, highlighting the importance of each stage. All experiments executed using a fixed amount of 5B tokens for the entire distillation process.

| MODEL TYPE | STAGES APPLIED | WINO.G. ACC \uparrow | ARC-E ACC \uparrow | ARC-C ACC \uparrow | PIQA ACC \uparrow | HELLAS. ACC \uparrow | LAMB. ACC \uparrow | AVG. ACC \uparrow |
|--------------|----------------|------------------------|----------------------|----------------------|---------------------|------------------------|----------------------|---------------------|
| Phi-Mamba | 2 | 55.9 | 75.4 | 38.0 | 75.2 | 56.6 | 18.9 | 53.3 |
| H. Phi-Mamba | 2 | 66.0 | 75.0 | 38.0 | 76.5 | 57.0 | 36.1 | 58.1 |
| Phi | 2 | 71.0 | 77.1 | 40.8 | 77.8 | 60.9 | 51.3 | 63.1 |
| Phi-Mamba | 3 | 62.8 | 64.3 | 27.8 | 75.6 | 52.6 | 43.8 | 54.5 |
| H. Phi-Mamba | 3 | 64.7 | 72.7 | 40.1 | 76.0 | 58.5 | 50.0 | 60.4 |
| Phi | 3 | 64.1 | 73.7 | 38.7 | 58.9 | 75.6 | 48.0 | 59.9 |
| Phi-Mamba | 2-3 | 72.3 | 75.0 | 40.8 | 75.2 | 59.7 | 50.6 | 62.3 |
| H. Phi-Mamba | 2-3 | 75.4 | 75.7 | 40.8 | 76.0 | 59.4 | 51.9 | 63.2 |
| Phi | 2-3 | 69.8 | 77.4 | 44.8 | 78.2 | 61.3 | 54.4 | 64.3 |
| Phi-Mamba | 1-3 | 74.8 | 72.7 | 43.5 | 75.6 | 59.6 | 49.2 | 62.7 |
| H. Phi-Mamba | 1-3 | 75.4 | 75.0 | 42.1 | 79.1 | 60.1 | 52.0 | 64.0 |
| Phi | 1-3 | 74.2 | 76.7 | 43.5 | 78.6 | 61.7 | 54.2 | 64.9 |

As elaborated in Section 5.7, this phase can freeze all network components except the Mamba-2 sequence mixer without a significant performance drop. This in particular indicates that the third stage (like the other stages of MOHAWK) can operate in computationally limited settings, enabling more users to utilize the MOHAWK distillation process.

In contrast to other phases of MOHAWK, we have observed occasional loss spikes during this phase. These abrupt spikes are typically seen in the training of large-scale language models and can negatively impact the model. We addressed this issue using checkpointing, weight decay, and gradient clipping, resulting in a more stable Stage 3.

5.3 Stage 2 (Hidden-State Alignment)

Following the analysis of the model’s end-to-end distillation in Stage 3, we evaluate the impact of aligning the hidden-state outputs of mixer blocks (Stage 2) on both the subsequent Stage 3 process and overall downstream model performance. We accomplish this by training Phi-Mamba instances from scratch using Stage 2 to various token counts. From these checkpoints, we proceed to Stage 3 training, ending with different total budgets to allow us to analyze how the degree of Stage 2 “pretraining” impacts Stage 3 performance at various token budgets.

Figure 3 demonstrates that given an adequate training budget, models beginning with weights with lower hidden state distances (after Stage 2) outperform those that depend exclusively on knowledge distillation (Stage 3). These lower hidden states are also correlated with lower starting perplexities, which in turn are correlated with downstream performance, as shown in Figure 5. Furthermore, Table 3 shows the synergy between Stage 2 and Stage 3, as applying Stage 3 on top of Stage

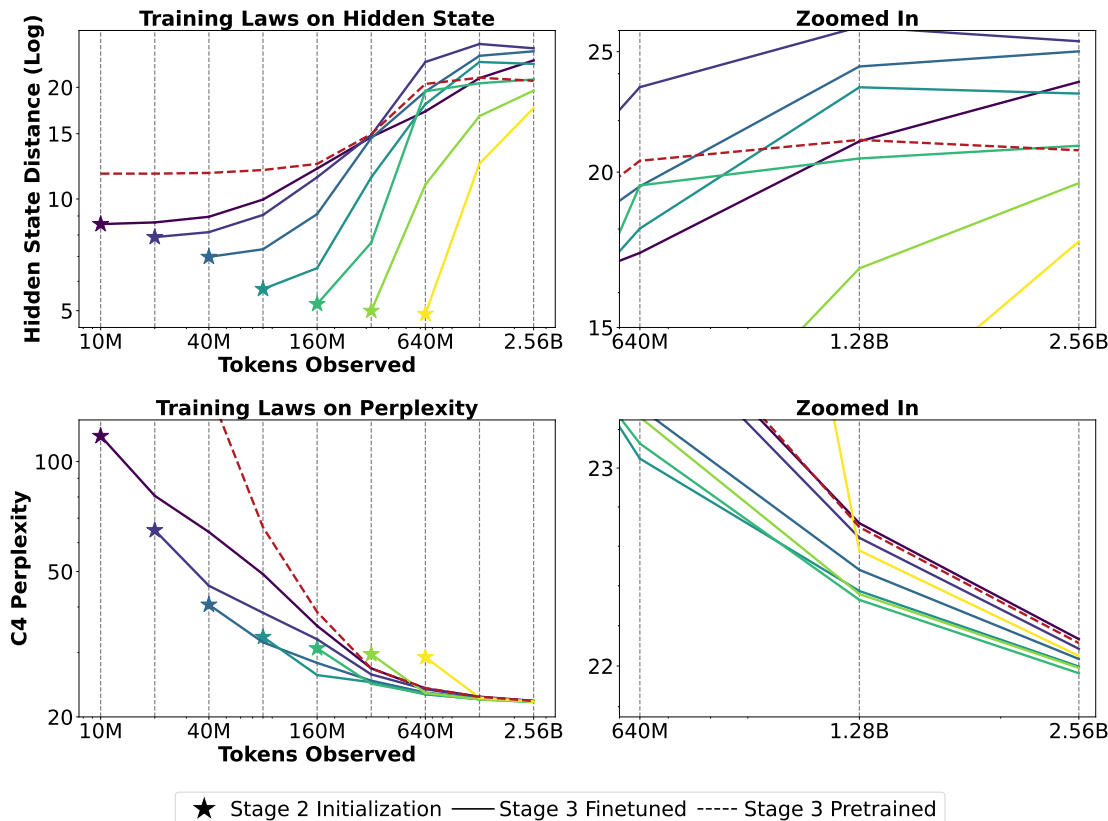


Figure 3: Training laws comparing the amount of token budget between Stages 2 and 3, as measured by the Stage 2 metric (hidden state distance) and Stage 3 metric (perplexity). Stage 2 initializations are used as the starting checkpoint for their respective Stage 3 finetuning models. Stage 3 pretrained is trained from scratch only with weight transfer and knowledge distillation. Despite training for less tokens on Stage 3 than the Stage 3 from scratch, almost all Stage 2 initialized models eventually outperform the baseline in perplexity on a fixed budget. In general, better aligned Stage 2 initializations improve post-Stage 3 performance.

2 outperforms vanilla knowledge distillation, highlighting the importance of incorporating both hidden-state alignment and knowledge distillation methods for the tested architectures. Table 3 also indicates that in scenarios where only this stage is applied, the closer the student architecture aligns with the teacher architecture (particularly, the more layers of attention it shares with Phi), the greater the impact of this stage on overall performance. However, when combining Stage 2 and 3, student models that are less similar to the teacher model have more noticeable improvements in their performance, e.g., the improvement for Phi-Mamba which has zero attention layers is larger than its hybrid counterpart which has four. We continue to explore the impact of Stage 2 on the downstream performance in Figure 5.

5.4 Stage 1 (Matrix Mixer Orientation)

Motivated by our previous finding, we then analyze how matching the matrix mixers can decrease the overall mixer block’s hidden-state distance with the teacher model even further. Similarly to our previous protocol, we assess the positive impact of the current stage on the following phase’s metrics and final model’s performance by comparing models with varying amount of Stage 1 and Stage 2 training on both stage metrics.

Figure 4 shows that even with constrained budgets, performing Stage 1 for a small period can help with subsequent stages and their performances. Thus, even a small amount of Stage 1 training can help their respective Stage 2 models reach better hidden-state distances compared to the from-scratch counterpart. This is despite the phenomenon that the teacher and student mixers diverge and then re-converge in Stage 2 after mixer similarity is no longer directly optimized. Coupled with Section 5.3, which discovers that lower hidden state initializations lead to better perplexity and downstream

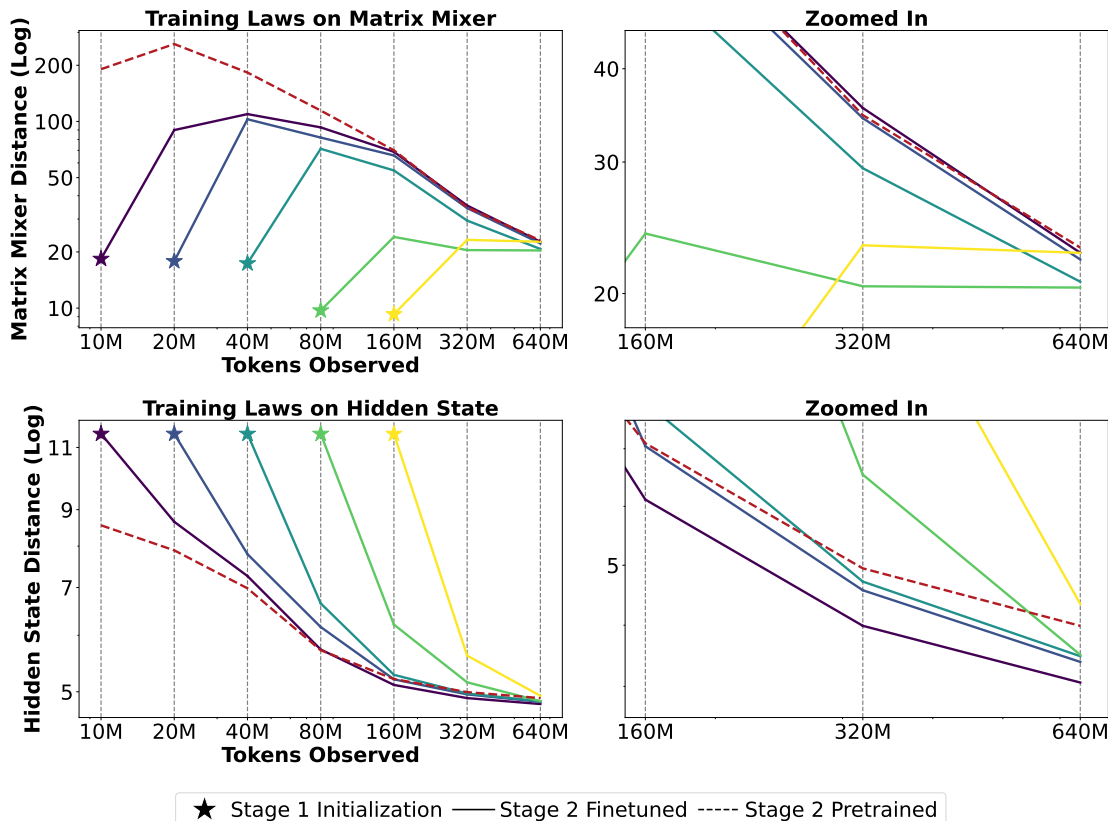


Figure 4: Training laws comparing the amount of token budget between Stages 1 and 2, as measured by the Stage 1 metric (matrix mixer distance) and Stage 2 metric (hidden state distance). Even a small amount of Stage 1 training can improve the model’s hidden-state distances in subsequent stages. Notably, this improvement occurs despite an increase in matrix mixer distance during Stage 2. This suggests that early Stage 1 training provides a foundational benefit that enhances the model’s performance in later stages, demonstrating the importance of initial training phases in model optimization.

performance, it can be inferred that Stage 1 aids the overall distillation process.

Furthermore, we empirically validate this intuition in Table 3, which indicates that this stage aligns the matrix mixers to a stronger degree than only the hidden state alignment. For example, employing only Stage 2 and 3 for Phi-to-Phi distillation does not allow the student to fully recover the original Phi-1.5’s performance on key metrics. Only by incorporating Stage 1 does the metric performance align with the original Phi teacher. Metric performance gains from the addition of Stage 1 can also be seen in both Phi-Mamba and Hybrid-Phi-Mamba.

5.5 Training the Final Phi-Mamba Model

After confirming the importance of the stages in Section 5.2, Section 5.3, and Section 5.4, we proceed to distill the final Phi-Mamba model using the three elements of MOHAWK. We use 80M tokens for Stage 1, due to the strong performance of the token count in both the matrix and hidden state distances (Figure 4). Stage 2 was distilled for 160M tokens given the apparent saturation of both hidden state distance and perplexity compared to the other initialization states, such as 10M, 20M, 40M, etc. (Figure 3). We employed Stage 3 to a total of 3B tokens across all stages and observed that the previously optimal learning rate applied for training training laws led to instabilities in training, particularly spikes in evaluation perplexity. Decreasing the learning rate for Stage 3 mitigated this issue (Appendix A). We hypothesize that the instability is due to the Stage 1 + 2 initialization’s Mamba component being quite similar to that of the teacher model, so a large learning rate coupled with disconnect between blocks, which are mended in Stage 3, can cause training instabilities. The performance of the final model is reported in Table 1.

Table 4: Given a budget to maintain four attention layers, we explore typical configurations to interleave them into the architecture. Phi-1.5 comprises 24 Attention layers, of which 20 are transformed into Mamba-2 blocks during MOHAWK, resulting in Hybrid-Phi-Mamba-1.5B, while the remaining 4 layers stay unchanged. The average block distance was taken at the start of the experiment.

| ATTENTION LAYERS | AVG BLOCK L2 DIST ↓ | WINO.G. Acc ↑ | ARC-E Acc ↑ | ARC-C Acc ↑ | PIQA Acc ↑ | HELLAS. Acc ↑ | LAMB. Acc ↑ | AVG. Acc ↑ |
|------------------|---------------------|---------------|-------------|-------------|------------|---------------|-------------|------------|
| 1-4 | 10.4 | 69.8 | 75.4 | 43.5 | 74.3 | 59.8 | 52.7 | 62.6 |
| 1,7,13,19 | 9.8 | 71.0 | 76.7 | 41.4 | 76.0 | 59.7 | 51.6 | 62.8 |
| 3,9,15,21 | 9.6 | 73.5 | 75.7 | 42.1 | 76.5 | 60.1 | 52.7 | 63.4 |
| 5,12,17,23 | 9.6 | 75.4 | 75.0 | 42.1 | 79.1 | 60.1 | 52.0 | 64.0 |
| 21-24 | 8.6 | 72.9 | 74.7 | 38.7 | 76.5 | 59.8 | 51.4 | 62.2 |

Table 5: Examination of how changing the number of interleaved Attention layers affects performance. In these experiments, we kept Attention layers 5, 12, 17, and 23 unchanged (see Table 4), and utilized MOHAWK to distill the intermediate layers. The average block distance was taken at the start of the experiment.

| # ATTENTION LAYERS | AVG BLOCK L2 DIST ↓ | WINO.G. Acc ↑ | ARC-E Acc ↑ | ARC-C Acc ↑ | PIQA Acc ↑ | HELLAS. Acc ↑ | LAMB. Acc ↑ | AVG. Acc ↑ |
|--------------------|---------------------|---------------|-------------|-------------|------------|---------------|-------------|------------|
| 1 | 11.0 | 69.8 | 72.3 | 41.4 | 76.5 | 60.4 | 50.2 | 61.8 |
| 2 | 10.2 | 72.9 | 72.7 | 40.8 | 76.5 | 60.0 | 51.0 | 62.3 |
| 4 | 8.6 | 75.4 | 75.0 | 42.1 | 79.1 | 60.1 | 52.0 | 64.0 |

5.6 Hybrid Phi-Mamba Model

Recently, models that integrate both Attention mechanisms and SSM layers have been proposed (Hatamizadeh and Kautz 2024; Lieber et al. 2024; Ren et al. 2024), delivering better results than using either architecture independently. Empirically, incorporating a limited number of Attention layers does make the training and inference time quadratic, although this effect is mitigated by the small number of Attention layers used.

We distill the Phi-1.5 model into a hybrid version, preserving only four original Attention layers and converting all other Attention blocks to Mamba-2 blocks through MOHAWK. This hybrid model achieves a downstream evaluation score of 66.0 (refer to Table 2), closely approaching the performance of the pure Attention Transformer architecture and exceeding the Phi-Mamba average score of 65.1. Hybrid-Phi-Mamba also performs well compared to other Attention-Mamba hybrids at the 1.5B size range while using less Attention layers and less overall parameters.

Table 4 shows the results for the most common placements of the four Attention layers in hybrid models. Despite all placements showing strong results, our experiments indicate that interleaving Mamba-2 layers uniformly yields superior performance on downstream evaluation benchmarks. This aligns with the solutions proposed by Samba (Ren et al. 2024), which also find that interleaving Attention layers within Mamba layers leads to improved performance.

Table 5 examines the impact of varying the number of interleaved Attention layers. Based on previous findings in Table 4, we carry out these experiments without converting the respective Attention layers in the network while utilizing MOHAWK to distill other layers. As anticipated, preserving a greater number of Attention layers results in improved outcomes. However, we hypothesize that there is still room for improvement for distilling hybrid models due to potential variations in the distillation process for hybrid versus non-hybrid architectures. These aspects, e.g., additional gradient updates, changes in optimizer settings, etc, could be further optimized, and we leave it for future work.

5.7 Approximating Self-Attention

Given the impact that Stage 1 (Matrix Orientation) and Stage 2 (Hidden-State Alignment) have on Stage 3’s (Weight Transfer and Knowledge-Distillation) effectiveness, we delve deeper into Mamba-2’s capability to learn interactions taught by Self-Attention. We first examine in Section 5.7.1 the extent to which a Mamba-2 sequence transformation can approximate a self-attention matrix. Next, we investigate in Section 5.7.2 whether this capability is evident in an end-to-end language model such as Phi-1.5.

Table 6: Attention matrix approximation by structured matrix mixers (Frobenius distance; lower is better). Structures are Toeplitz, (causal) low-rank (LR), state space dual (SSD) model (3.2) and general semi-separable matrices (SSM). We have used 1,000 samples, each consisting of 512 tokens. Llama2-7B-Chat was applied on every sample, and one attention head from each layer was randomly chosen for approximation. We evaluated (LR) and SSD families with 10,000 gradient descent steps per sample.

| STRUCTURE (State size N) | TOEP. | LR (16) | SSD (16) | SSM (16) | LR (32) | SSD (32) | SSM (32) | LR (64) | SSD (64) | SSM (64) |
|--------------------------------|-------|------------|-------------|-------------|------------|-------------|-------------|------------|-------------|-------------|
| WT-103 | 12.0 | 0.619 | 0.477 | 0.266 | 0.322 | 0.237 | 0.127 | 0.132 | 0.097 | 0.046 |
| OWT | 12.2 | 0.606 | 0.466 | 0.259 | 0.314 | 0.231 | 0.123 | 0.129 | 0.095 | 0.045 |
| C4 | 12.3 | 0.595 | 0.453 | 0.236 | 0.310 | 0.226 | 0.112 | 0.128 | 0.093 | 0.041 |
| IMdB | 12.3 | 0.598 | 0.455 | 0.238 | 0.312 | 0.226 | 0.113 | 0.129 | 0.094 | 0.043 |

5.7.1 Self-Attention Approximation with Structured Matrix Mixers

We start by testing the ability of various matrix mixer families to match the empirical self-attention matrices of a pretrained Transformer. We take 1000 samples from each layer of a Llama2-7b-Chat model (Touvron et al. 2023), materialize the attention matrices, and project them onto given classes of structured matrices. The results in Table 6 are averaged across all layers.

In particular, to describe the class of linear attention matrices (3.1), we use the fact that \mathbf{Q} and \mathbf{K} are projections of the input $x \in \mathbb{R}^{d_{in}}$ onto $\mathbb{R}^{d_{out}}$, and therefore their rank is bounded by $\min\{d_{in}, d_{out}\}$. For multihead linear attention, d_{out} (also known as head dimension) is typically a small value (e.g., Phi-1.5 and Llama2-7b-Chat have head dimensions of 64 and 128, respectively). Thus, we approximate this family of sequence mixers using causal low-rank matrices $\mathbf{L} \circ \mathbf{Q}\mathbf{K}^T$, where \mathbf{L} is a lower-triangular causal mask of 1s, and \mathbf{Q}, \mathbf{K} are in $\mathbb{R}^{n \times d}$ with $d \ll n$ (indicating that the head dimension is substantially smaller than the sequence length).

To describe the multi-head Mamba-2 matrix family, we utilize the state space dual (SSD) layer (3.2) in a manner similar to the previous linear attention, but now the causal matrix \mathbf{L} possesses an n -degree rolling multiplicative structure for SSD which can be seen as a more expressive mask that generalizes the causal mask (Section 3.2).

Both causal low-rank and SSD matrix families were approximated with 10,000 steps of gradient descent per sample. Finally, to approximate the general class of SSM matrix mixers, we utilize *balanced truncation*, a projection algorithm that does not rely on gradients. This method is mainly known in the field of time-invariant Dynamical System model reduction (Gugercin and Antoulas 2004) and has been modified for use in time-varying systems (Sandberg and Rantzer 2004). Similarly, for the family of causal Toeplitz matrices, which represent a convolution operation, we employ a simple heuristic that minimizes the error for each attention matrix.

Table 6 shows that while the SSM matrix family provides the closest approximation to the self-attention matrix mixer, the Mamba-2 mixer family (SSD) has just twice the distance from the SSM matrices. This is in contrast to Linear Attention, which has three times the distance, all while keeping a computational cost on par with Linear Attention. More details can be found in Appendix C.

5.7.2 Self-Attention Replacement in Language Models

Since the experiments in Section 5.7.1 were designed to approximate a self-attention matrix under controlled conditions, we further validate the ability of a Mamba-2 block to replace an Attention layer within a language model. Firstly, we create two variants of our architecture, Phi-Toeplitz and Phi-LR, and run the MOHAWK process for 1B tokens at each stage (see Table 7) to verify that the previous finding hold in a multilayer, end-to-end model case.

Secondly, we run MOHAWK while freezing various parts of the Phi-Mamba modules (refer to Table 8), revealing that limiting the trainable elements to the Mamba-2 blocks (excluding the embedding, head and all MLP layers) results in only a minor performance decrease during MOHAWK distillation.

Interestingly, in all of the aforementioned experiments, we have found a consistent correlation between the projection distances of the matrix (Frobenius distance) in Table 6 and the downstream performance metrics (accuracy) in Table 7. Essentially, a better matrix approximation (lower Frobenius distance) is correlated with better model performance (higher

Table 7: Ablations of matrix structure using the same training recipe (Stages 2 and 3). While many efficient sequence models (e.g. global convolutions, linear attention, and state space models) can be represented as structured matrix mixers (e.g. Toeplitz, low-rank, and semi-separable matrices respectively), more expressive structured matrix families can match the attention matrix more closely.

| MATRIX STRUCTURE | BLOCK OUTPUT L2 DIST. ↓ | WINO.G. Acc ↑ | ARC-E Acc ↑ | ARC-C Acc ↑ | PIQA Acc ↑ | HELLAS. Acc ↑ | AVG. Acc ↑ |
|------------------|-------------------------|---------------|-------------|-------------|-------------|---------------|-------------|
| Causal Toeplitz | 9.6 | 49.3 | 21.2 | 26.2 | 52.3 | 25.9 | 35.0 |
| Causal low-rank | 7.6 | 50.2 | 27.9 | 25.0 | 53.3 | 25.6 | 36.4 |
| SSD | 5.5 | 67.2 | 71.0 | 38.6 | 74.2 | 45.0 | 59.2 |

accuracy) on various tasks. This connection highlights the relationship between the quality of the matrix approximation and the performance of the model. Such findings are echoed in Hwang et al. (2024), which find that more expressive matrix mixers lead to more performant models, e.g., Low-rank-based BERT models outperform Toeplitz-based ones.

Table 8: Distillation with MOHAWK for both Phi-Mamba-1.5B and Hybrid-Phi-Mamba-1.5B (with the final four Attention layers unchanged). MOHAWK can be employed while maintaining all components other than the sequence mixer blocks frozen without compromising Phi-Mamba’s performance (Section 5.2).

| MODEL TYPE | TRAINABLE COMPONENTS | WINO.G. Acc ↑ | ARC-E Acc ↑ | ARC-C Acc ↑ | PIQA Acc ↑ | HELLAS. Acc ↑ | LAMB. Acc ↑ | AVG. Acc ↑ |
|--------------|----------------------|---------------|-------------|-------------|------------|---------------|-------------|-------------|
| Phi-Mamba | All | 74.8 | 72.7 | 43.5 | 75.6 | 59.6 | 49.2 | 62.7 |
| | Mamba-2 | 69.1 | 73.5 | 43.8 | 74.7 | 59.3 | 48.2 | 61.4 |
| Hybrid Phi.M | All | 75.4 | 75.0 | 42.1 | 79.1 | 60.1 | 52.0 | 64.0 |
| | Mamba-2 | 78.6 | 75.0 | 41.4 | 76.5 | 59.8 | 51.9 | 63.9 |

6 Discussion and Conclusion

Our experiments shows that the Mamba-2 model can be successfully distilled from a pretrained Transformer teacher model, utilizing its extensive knowledge learned from custom datasets and higher computational resources. Despite using less than 100× data compared to many open-source models, including Mamba, our subquadratic model outperforms other subquadratic models in various benchmark tests by a wide margin.

The MOHAWK framework’s multi-stage process which gradually increased the scope of distillation is essential extracting the teacher model’s knowledge to the fullest extent as shown in our ablations and training laws. We continue to find the effectiveness of MOHAWK when distilling hybrid Attention-SSM models and provide ablations on the number and position of Attention layers.

Additionally, we demonstrate that Mamba-2’s relationship to Transformers is evident not only in theory, but also in practice, as it captures interactions similar to those of Transformers, and is able to replace Attention with little drop in performance. Coupled with past research which has posited that much of a language model’s knowledge is embedded in the MLP blocks, we believe that any subquadratic model with a sufficiently expressive matrix mixer can replicate the behavior of pretrained Transformers, bringing quadratic knowledge to subquadratic models. We recommend further research to explore the role of sequence mixing layers in subquadratic models and their impact on performance. Advancements in both the distillation process and the sequence mixer architecture could lead to further improved performance in a range of tasks. We propose that “trainability” and “distillability” are distinct properties of the models, and therefore, distillation techniques should be more appropriately tailored to the model.

References

- [1] Maximilian Beck, Korbinian Pöppel, Markus Spanring, Andreas Auer, Oleksandra Prudnikova, Michael Kopp, Günter Klambauer, Johannes Brandstetter, and Sepp Hochreiter. *xLSTM: Extended Long Short-Term Memory*. 2024. arXiv: [2405.04517 \[cs.LG\]](#).
- [2] Yonatan Bisk, Rowan Zellers, Jianfeng Gao, Yejin Choi, et al. “PIQA: Reasoning about Physical Commonsense in Natural Language”. In: *Proceedings of the AAAI conference on Artificial Intelligence*. Vol. 34. 05. 2020, pp. 7432–7439.
- [3] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. *Language Models are Few-Shot Learners*. 2020. arXiv: [2005.14165 \[cs.CL\]](#).
- [4] Shiv Chandrasekaran, Patrick Dewilde, Ming Gu, T Pals, and Alle-Jan van der Veen. “Fast stable solver for sequentially semi-separable linear systems of equations”. In: *International Conference on High-Performance Computing*. Springer. 2002, pp. 545–554.
- [5] Yen-Chun Chen, Zhe Gan, Yu Cheng, Jingzhou Liu, and Jingjing Liu. *Distilling Knowledge Learned in BERT for Text Generation*. 2020. arXiv: [1911.03829 \[cs.CL\]](#).
- [6] 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”. In: *arXiv preprint arXiv:1803.05457* (2018).
- [7] Tri Dao, Daniel Y. Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. *FlashAttention: Fast and Memory-Efficient Exact Attention with IO-Awareness*. 2022. arXiv: [2205.14135 \[cs.LG\]](#).
- [8] Tri Dao and Albert Gu. “Transformers are SSMS: Generalized Models and Efficient Algorithms Through Structured State Space Duality”. In: *International Conference on Machine Learning (ICML)*. 2024.
- [9] P. Dewilde and A. J. van der Veen. “On the Hankel-norm approximation of upper-triangular operators and matrices”. In: *Integral Equations and Operator Theory* 17.1 (Mar. 1, 1993), pp. 1–45. DOI: [10.1007/BF01322544](#). URL: <https://doi.org/10.1007/BF01322544>.
- [10] Patrick Dewilde and Alle-Jan Veen. *Time-Varying Systems and Computations*. 1st ed. Springer Book Archive. Springer Science+Business Media Dordrecht 1998. Springer New York, NY, 1998, pp. XIV, 460. ISBN: 978-1-4757-2817-0. DOI: <https://doi.org/10.1007/978-1-4757-2817-0>. URL: <https://doi.org/10.1007/978-1-4757-2817-0>.
- [11] Daniel Y. Fu, Tri Dao, Khaled K. Saab, Armin W. Thomas, Atri Rudra, and Christopher Ré. *Hungry Hungry Hippos: Towards Language Modeling with State Space Models*. 2023. arXiv: [2212.14052 \[cs.LG\]](#).
- [12] Albert Gu. “Modeling Sequences with Structured State Spaces”. PhD thesis. Stanford University, 2023.
- [13] Albert Gu and Tri Dao. *Mamba: Linear-Time Sequence Modeling with Selective State Spaces*. 2023. arXiv: [2312.00752 \[cs.LG\]](#).
- [14] Albert Gu, Karan Goel, and Christopher Ré. *Efficiently Modeling Long Sequences with Structured State Spaces*. 2022. arXiv: [2111.00396 \[cs.LG\]](#).
- [15] Serkan Gugercin and Athanasios C. Antoulas. “A Survey of Model Reduction by Balanced Truncation and Some New Results”. In: *International Journal of Control* 77.8 (2004), pp. 748–766. DOI: [10.1080/00207170410001713448](#). eprint: <https://doi.org/10.1080/00207170410001713448>. URL: <https://doi.org/10.1080/00207170410001713448>.
- [16] Suriya Gunasekar, Yi Zhang, Jyoti Aneja, Caio César Teodoro Mendes, Allie Del Giorno, Sivakanth Gopi, Mojan Javaheripi, Piero Kauffmann, Gustavo de Rosa, Olli Saarikivi, Adil Salim, Shital Shah, Harkirat Singh Behl, Xin Wang, Sébastien Bubeck, Ronen Eldan, Adam Tauman Kalai, Yin Tat Lee, and Yuanzhi Li. *Textbooks Are All You Need*. 2023. arXiv: [2306.11644 \[cs.CL\]](#).
- [17] Sangchul Hahn and Heeyoul Choi. *Self-Knowledge Distillation in Natural Language Processing*. 2019. arXiv: [1908.01851 \[cs.CL\]](#).
- [18] Md. Akmal Haidar and Mehdi Rezagholizadeh. *TextKD-GAN: Text Generation using KnowledgeDistillation and Generative Adversarial Networks*. 2019. arXiv: [1905.01976 \[cs.CL\]](#).

- [19] Ali Hatamizadeh and Jan Kautz. *MambaVision: A Hybrid Mamba-Transformer Vision Backbone*. 2024. arXiv: [2407.08083](https://arxiv.org/abs/2407.08083) [cs.CV]. URL: <https://arxiv.org/abs/2407.08083>.
- [20] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. *Distilling the Knowledge in a Neural Network*. 2015. arXiv: [1503.02531](https://arxiv.org/abs/1503.02531) [stat.ML].
- [21] Minghao Hu, Yuxing Peng, Furu Wei, Zhen Huang, Dongsheng Li, Nan Yang, and Ming Zhou. *Attention-Guided Answer Distillation for Machine Reading Comprehension*. 2018. arXiv: [1808.07644](https://arxiv.org/abs/1808.07644) [cs.CL].
- [22] Sukjun Hwang, Aakash Lahoti, Tri Dao, and Albert Gu. *Hydra: Bidirectional State Space Models Through Generalized Matrix Mixers*. 2024. arXiv: [2407.09941](https://arxiv.org/abs/2407.09941) [cs.LG]. URL: <https://arxiv.org/abs/2407.09941>.
- [23] Ananya Harsh Jha, Dirk Groeneveld, Emma Strubell, and Iz Beltagy. “Large Language Model Distillation Doesn’t Need a Teacher”. In: *arXiv preprint arXiv:2305.14864* (2023).
- [24] Jungo Kasai, Hao Peng, Yizhe Zhang, Dani Yogatama, Gabriel Ilharco, Nikolaos Pappas, Yi Mao, Weizhu Chen, and Noah A Smith. “Finetuning pretrained transformers into rnns”. In: *arXiv preprint arXiv:2103.13076* (2021).
- [25] Angelos Katharopoulos, Apoorv Vyas, Nikolaos Pappas, and François Fleuret. *Transformers are RNNs: Fast Autoregressive Transformers with Linear Attention*. 2020. arXiv: [2006.16236](https://arxiv.org/abs/2006.16236) [cs.LG].
- [26] Chen Liang, Haoming Jiang, Zheng Li, Xianfeng Tang, Bin Yin, and Tuo Zhao. “Homodistil: Homotopic task-agnostic distillation of pre-trained transformers”. In: *arXiv preprint arXiv:2302.09632* (2023).
- [27] Opher Lieber, Barak Lenz, Hofit Bata, Gal Cohen, Jhonathan Osin, Itay Dalmedigos, Erez Safahi, Shaked Meirom, Yonatan Belinkov, Shai Shalev-Shwartz, Omri Abend, Raz Alon, Tomer Asida, Amir Bergman, Roman Glozman, Michael Gokhman, Avshalom Manevich, Nir Ratner, Noam Rozen, Erez Shwartz, Mor Zusman, and Yoav Shoham. *Jamba: A Hybrid Transformer-Mamba Language Model*. 2024. arXiv: [2403.19887](https://arxiv.org/abs/2403.19887) [cs.CL]. URL: <https://arxiv.org/abs/2403.19887>.
- [28] Hao Liu, Matei Zaharia, and Pieter Abbeel. *Ring Attention with Blockwise Transformers for Near-Infinite Context*. 2023. arXiv: [2310.01889](https://arxiv.org/abs/2310.01889) [cs.CL].
- [29] Harsh Mehta, Ankit Gupta, Ashok Cutkosky, and Behnam Neyshabur. *Long Range Language Modeling via Gated State Spaces*. 2022. arXiv: [2206.13947](https://arxiv.org/abs/2206.13947) [cs.LG].
- [30] Samuel A. Melchior, Paul Van Dooren, and Kyle A. Gallivan. “Model reduction of linear time-varying systems over finite horizons”. In: *Applied Numerical Mathematics* 77 (2014), pp. 72–81. ISSN: 0168-9274. DOI: <https://doi.org/10.1016/j.apnum.2013.10.007>. URL: <https://www.sciencedirect.com/science/article/pii/S0168927413001414>.
- [31] Jean Mercat, Igor Vasiljevic, Sedrick Keh, Kushal Arora, Achal Dave, Adrien Gaidon, and Thomas Kollar. *Linearizing Large Language Models*. 2024. arXiv: [2405.06640](https://arxiv.org/abs/2405.06640) [cs.CL].
- [32] Jingcheng Niu, Andrew Liu, Zining Zhu, and Gerald Penn. *What does the Knowledge Neuron Thesis Have to do with Knowledge?* 2024. arXiv: [2405.02421](https://arxiv.org/abs/2405.02421) [cs.CL].
- [33] Denis Paperno, Germán Kruszewski, Angeliki Lazaridou, Ngoc-Quan Pham, Raffaella Bernardi, Sandro Pezzelle, Marco Baroni, Gemma Boleda, and Raquel Fernández. “The LAMBADA Dataset: Word Prediction Requiring a Broad Discourse Context”. In: *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics*. 2016, pp. 1525–1534.
- [34] Bo Peng, Eric Alcaide, Quentin Anthony, Alon Albalak, Samuel Arcadinho, Stella Biderman, Huanqi Cao, Xin Cheng, Michael Chung, Matteo Grella, Kranthi Kiran GV, Xuzheng He, Haowen Hou, Jiaju Lin, Przemyslaw Kazienko, Jan Kocon, Jiaming Kong, Bartłomiej Koptyra, Hayden Lau, Krishna Sri Ipsit Mantri, Ferdinand Mom, Atsushi Saito, Guangyu Song, Xiangru Tang, Bolun Wang, Johan S. Wind, Stanislaw Wozniak, Ruichong Zhang, Zhenyuan Zhang, Qihang Zhao, Peng Zhou, Qinghua Zhou, Jian Zhu, and Rui-Jie Zhu. *RWKV: Reinventing RNNs for the Transformer Era*. 2023. arXiv: [2305.13048](https://arxiv.org/abs/2305.13048) [cs.CL].
- [35] Zhen Qin, Xiaodong Han, Weixuan Sun, Bowen He, Dong Li, Dongxu Li, Yuchao Dai, Lingpeng Kong, and Yiran Zhong. *Toeplitz Neural Network for Sequence Modeling*. 2023. arXiv: [2305.04749](https://arxiv.org/abs/2305.04749) [cs.CL].
- [36] Zhen Qin, Songlin Yang, Weixuan Sun, Xuyang Shen, Dong Li, Weigao Sun, and Yiran Zhong. *HGRN2: Gated Linear RNNs with State Expansion*. 2024. arXiv: [2404.07904](https://arxiv.org/abs/2404.07904) [cs.CL].

- [37] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. *Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer*. 2023. arXiv: [1910.10683](https://arxiv.org/abs/1910.10683) [cs.LG]. URL: <https://arxiv.org/abs/1910.10683>.
- [38] Liliang Ren, Yang Liu, Yadong Lu, Yelong Shen, Chen Liang, and Weizhu Chen. *Samba: Simple Hybrid State Space Models for Efficient Unlimited Context Language Modeling*. 2024. arXiv: [2406.07522](https://arxiv.org/abs/2406.07522) [cs.CL]. URL: <https://arxiv.org/abs/2406.07522>.
- [39] Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. “Winogrande: An Adversarial Winograd Schema Challenge at Scale”. In: *Communications of the ACM* 64.9 (2021), pp. 99–106.
- [40] H. Sandberg and A. Rantzer. “Balanced truncation of linear time-varying systems”. In: *IEEE Transactions on Automatic Control* 49.2 (2004), pp. 217–229. DOI: [10.1109/TAC.2003.822862](https://doi.org/10.1109/TAC.2003.822862).
- [41] 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. arXiv: [2307.08621](https://arxiv.org/abs/2307.08621) [cs.CL].
- [42] Xu Tan, Yi Ren, Di He, Tao Qin, Zhou Zhao, and Tie-Yan Liu. *Multilingual Neural Machine Translation with Knowledge Distillation*. 2019. arXiv: [1902.10461](https://arxiv.org/abs/1902.10461) [cs.CL].
- [43] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. *LLaMA: Open and Efficient Foundation Language Models*. 2023. arXiv: [2302.13971](https://arxiv.org/abs/2302.13971) [cs.CL].
- [44] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. *Attention Is All You Need*. 2023. arXiv: [1706.03762](https://arxiv.org/abs/1706.03762) [cs.CL].
- [45] A.J. van der Veen and P. Dewilde. “On low-complexity approximation of matrices”. In: *Linear Algebra and its Applications* 205-206 (1994), pp. 1145–1201. ISSN: 0024-3795. DOI: [https://doi.org/10.1016/0024-3795\(94\)90383-2](https://doi.org/10.1016/0024-3795(94)90383-2). URL: <https://www.sciencedirect.com/science/article/pii/0024379594903832>.
- [46] Jue Wang, Wentao Zhu, Pichao Wang, Xiang Yu, Linda Liu, Mohamed Omar, and Raffay Hamid. *Selective Structured State-Spaces for Long-Form Video Understanding*. 2023. arXiv: [2303.14526](https://arxiv.org/abs/2303.14526) [cs.CV].
- [47] Wenhui Wang, Furu Wei, Li Dong, Hangbo Bao, Nan Yang, and Ming Zhou. “Minilm: Deep self-attention distillation for task-agnostic compression of pre-trained transformers”. In: *Advances in Neural Information Processing Systems* 33 (2020), pp. 5776–5788.
- [48] Zicheng Wang, Zhenghao Chen, Yiming Wu, Zhen Zhao, Luping Zhou, and Dong Xu. *PoinTramba: A Hybrid Transformer-Mamba Framework for Point Cloud Analysis*. 2024. arXiv: [2405.15463](https://arxiv.org/abs/2405.15463) [cs.CV]. URL: <https://arxiv.org/abs/2405.15463>.
- [49] Mengzhou Xia, Tianyu Gao, Zhiyuan Zeng, and Danqi Chen. “Sheared llama: Accelerating language model pre-training via structured pruning”. In: *arXiv preprint arXiv:2310.06694* (2023).
- [50] Songlin Yang, Bailin Wang, Yikang Shen, Rameswar Panda, and Yoon Kim. *Gated Linear Attention Transformers with Hardware-Efficient Training*. 2024. arXiv: [2312.06635](https://arxiv.org/abs/2312.06635) [cs.LG].
- [51] Ze Yang, Linjun Shou, Ming Gong, Wutao Lin, and Daxin Jiang. *Model Compression with Two-stage Multi-teacher Knowledge Distillation for Web Question Answering System*. 2019. arXiv: [1910.08381](https://arxiv.org/abs/1910.08381) [cs.CL].
- [52] Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. “HellaSwag: Can a Machine Really Finish Your Sentence?” In: *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. 2019.
- [53] Michael Zhang, Kush Bhatia, Hermann Kumbong, and Christopher Ré. *The Hedgehog & the Porcupine: Expressive Linear Attentions with Softmax Mimicry*. 2024. arXiv: [2402.04347](https://arxiv.org/abs/2402.04347) [cs.LG].
- [54] Chunting Zhou, Graham Neubig, and Jiatao Gu. *Understanding Knowledge Distillation in Non-autoregressive Machine Translation*. 2021. arXiv: [1911.02727](https://arxiv.org/abs/1911.02727) [cs.CL].

A Experiments and Experimental Details

To construct Section 5.5, we performed grid searches for training in Stages 1, 2, and 3 independently from scratch to find the optimal hyperparameters. We explored learning rates $lr = \{1, 2, 5\} \times 10^{\{-3, -4\}}$ and batch sizes $2^{\{15, 16, 17, 18\}}$. AdamW Optimizer was used with $\beta = (0.9, 0.95)$, incorporating a weight decay of 0.1, gradient clipping at 1.0, and a Warmup-Stable-Decay (WSD) scheduler with 10% warmup and 10% decay utilizing linear warmup and cooldown functions. Automatic mixed precision training to bf16 was used in all stages. For Stages 1 and 2, we initially fixed the batch size at 2^{16} , then varied the learning rates. After identifying the optimal learning rate, we adjusted the batch sizes and subsequently finalized the learning rate after fixing the batch size. Consequently, Stage 1 used $bs = 2^{15}$, $lr = 5 \times 10^{-4}$ and Stage 2 used $bs = 2^{15}$, $lr = 2 \times 10^{-3}$. In Stage 3, we set the batch size to $2^{19} \approx 0.5M$ and focused solely on varying the learning rate, resulting in 5×10^{-4} . Stages 1 and 2 were trained to 200M steps each while Stage 3 extended to 1B steps. For the Phi-Mamba ultimate model, the Stage 3 learning rate was reduced to 2×10^{-4} to enhance stability.

In the development of the training law (see Figure 3), we executed a single "continuous" run initialized from a state that included several checkpoints. The warm-up period was determined as 10% of the tokens processed during the continuous run. For instance, if the model's goal was to process 640 million tokens, and it started from a run that had processed 40 million tokens, then the warm-up would be set at 60 million tokens. The checkpoints recorded during the warm-up phase were preserved as they were, while subsequent checkpoints underwent a cooling of 10% of the current phase. To illustrate, in the scenario mentioned earlier, a checkpoint at 320 million tokens during the 40M to 640M run would maintain the original warmup, while the cooldown would span 28 million tokens. Conversely, a checkpoint at 80 million tokens within the warm-up phase would be saved without any cooldown.

Figure 5 extends the Stage 2 versus Stage 3 comparison in Figure 3, except we measure average accuracy on downstream metrics instead of perplexity. We observe a strong correlation between the training laws of perplexity and downstream evaluation metrics. While the general trend indicates that models exposed to more tokens during the prior stage initialization tend to perform better on both perplexity and downstream metrics, the relationship is not perfectly aligned. Specifically, the order of model performance based on perplexity does not always match the order based on downstream metrics, highlighting some differences in how these metrics capture model effectiveness.

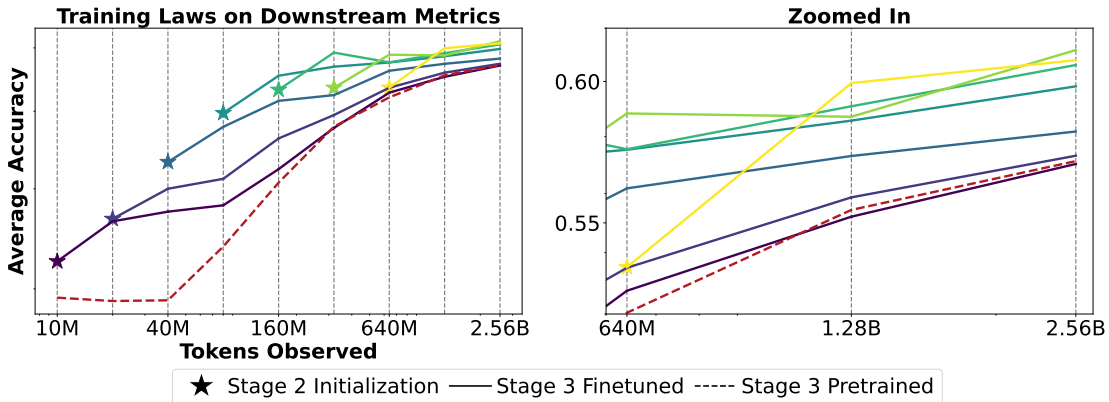


Figure 5: Training laws comparing the amount of token budget between Stages 2 and 3, as measured by the average accuracy of downstream evaluation metrics.

B Applying Mamba-2 as a Black Box

As noted previously Section 4.4, our Mamba-based sequence mixer is slightly modified from the original to make it more amenable for distilling from a Transformer architecture. In particular, the Mamba-2 sequence mixer is treated entirely in discrete time by projecting the input onto the matrix A and removing the discretization parameter Δ . Even though this formulation is somewhat different from Mamba-2, the original algorithm remains applicable through a reduction expressed in Appendix B.

Listing 1 PyTorch example for using the Mamba algorithm for a *Delta*-free variation.

```

"""
X: (batch, seqlen, nheads, headdim)
A_log: (batch, seqlen, nheads)
B: (batch, seqlen, nheads, dstate)
C: (batch, seqlen, nheads, dstate)
D: (nheads)
"""

y = Mamba(
    X = X / A_log.unsqueeze(-1),
    dt = rearrange(A_log, "b c h -> b h c"),
    A = torch.ones(self.nheads),
    B = B,
    C = C,
)
Du = torch.einsum("h,blhp->blhp", D, X)
y = rearrange(y + Du, "b l h p -> b l (h p)")

```

C Attention Matrix Approximation Details

This section serves as a complement to Section 5.7.1 and outlines the methods employed to create Table 6. Appendices C.1 to C.5 describe our strategies for finding a matrix within the specified families that closely approximates the original attention matrix using a selected distance metric. Formally, we consider the following optimization problem:

$$\min_{\mathbf{X} \in \mathcal{M}} \|\mathbf{M} - \mathbf{X}\| \tag{6}$$

where \mathcal{M} is the subspace of a specific matrix family, \mathbf{M} is the attention matrix, and $\|\cdot\|$ corresponds to a selected distance metric. In the following sections, we explore different methods and matrix families for this optimization problem.

C.1 Semi-Separable Matrix Approximation

Considering a time-varying system denoted by $\{\mathbf{A}_k, \mathbf{B}_k, \mathbf{C}_k, \mathbf{D}_k\}_{k \in [l]}$, we can describe it using the matrix mixer T (also known as the transfer matrix) as follows:

$$T = \begin{bmatrix} D_1 & 0 & 0 & 0 & 0 & \dots & 0 \\ C_2 B_1 & D_2 & 0 & 0 & 0 & \dots & 0 \\ C_3 A_2 B_1 & C_3 B_2 & D_3 & 0 & 0 & \dots & 0 \\ C_4 A_3 B_1 & C_4 A_3 B_2 & C_4 B_3 & D_4 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ C_l A_{l-1} B_1 & C_l A_{l-1} B_2 & C_l A_{l-1} B_3 & C_l A_{l-1} B_4 & \dots & C_l B_{l-1} & D_l \end{bmatrix}$$

With $\mathbf{A}_k \in \mathbb{R}^{n \times n}$, $\mathbf{B}_k \in \mathbb{R}^{m \times n}$, $\mathbf{C}_k \in \mathbb{R}^{p \times n}$, and $\mathbf{D}_k \in \mathbb{R}^{p \times m}$, where n is the state dimension, m the input dimension, and p the output dimension.

As T 's form corresponds to a semi-separable matrix (i.e., each sub-matrix has a rank of up to n), we will label this matrix form as *SSM with state size n* throughout the remainder of this appendix, representing a state-space model with state size n or a semi-separable matrix of order n .

Every matrix can be represented as a linear combination of rank-one matrices. Thus, the attention matrix $\mathbf{M} \in \mathbb{R}^{l \times l}$ can be interpreted as an SSM with a state size of up to $l \gg n$. Consequently, we can employ prior research on time-varying model order reduction (P. Dewilde and A. J. v. d. Veen 1993; Melchior, Van Dooren, and Gallivan 2014; A. v. d. Veen and P. Dewilde 1994) to reduce \mathbf{M} to an SSM with a smaller state size n . Specifically, we utilize the following SVD-based approximation:

Table 9: Full attention matrix approximation by structured matrix mixers Structures are Toeplitz, causal low-rank (LR), RetNet, state space dual (SSD) model (3.2) with and without the diagonal \mathbf{D} term and general semi-separable matrices (SSM). We have used 1,000 samples, each consisting of 512 tokens. Llama2-7B-Chat was applied on every sample, and one attention head from each layer was randomly chosen for approximation. We evaluated (LR), RetNet, and SSD families with 10,000 gradient descent steps per sample.

| STRUCTURE (State size N) | TOEPLITZ - | CAUSAL LOW-RANK (16) | RETNET (16) | SSD WITHOUT \mathbf{D} (16) | SSD (16) | SEMI. SEP. MATRIX (16) |
|--------------------------------|---------------|-------------------------|----------------|----------------------------------|-------------|---------------------------|
| WT-103 | 12.0 | 0.619 | 0.530 | 0.522 | 0.477 | 0.266 |
| OWT | 12.2 | 0.606 | 0.516 | 0.508 | 0.466 | 0.259 |
| C4 | 12.3 | 0.595 | 0.503 | 0.496 | 0.453 | 0.236 |
| IMdB | 12.3 | 0.598 | 0.505 | 0.498 | 0.455 | 0.238 |
| (State size N) | - | (32) | (32) | (32) | (32) | (32) |
| WT-103 | 12.0 | 0.322 | 0.268 | 0.262 | 0.237 | 0.127 |
| OWT | 12.2 | 0.314 | 0.261 | 0.255 | 0.231 | 0.123 |
| C4 | 12.3 | 0.310 | 0.255 | 0.249 | 0.226 | 0.112 |
| IMdB | 12.3 | 0.312 | 0.256 | 0.251 | 0.226 | 0.113 |
| (State size N) | - | (64) | (64) | (64) | (64) | (64) |
| WT-103 | 12.0 | 0.132 | 0.110 | 0.107 | 0.097 | 0.046 |
| OWT | 12.2 | 0.129 | 0.107 | 0.104 | 0.095 | 0.045 |
| C4 | 12.3 | 0.128 | 0.106 | 0.102 | 0.093 | 0.041 |
| IMdB | 12.3 | 0.129 | 0.106 | 0.103 | 0.094 | 0.043 |

Algorithm 1 Approximation of Attention Matrix \mathbf{M} as an SSM with State Size n

Input: Attention matrix $\mathbf{M} \in \mathbb{R}^{L \times L}$, state size n

Output: Approximated attention matrix $\tilde{\mathbf{M}} \in \mathbb{R}^{L \times L}$

Procedure:

1. For $k = 1, \dots, L - 1$:

1.1 Define H_k as the submatrix of \mathbf{M} below and to the left of entry $M_{k,k}$:

$$H_k = \begin{bmatrix} M_{k,1} & \cdots & M_{k,k-1} \\ \vdots & \ddots & \vdots \\ M_{L,1} & \cdots & M_{L,k-1} \end{bmatrix}$$

1.2 Perform the SVD on H_k and truncate it to rank n

1.3 Integrate the truncated H_k back into the new matrix $\tilde{\mathbf{M}}$

Note that the diagonal elements of \mathbf{M} were not subject to approximation in Algorithm 1 as they remain unchanged.

Although this approximation method for the semi-separable matrix is heuristic, it has been empirically shown to deliver good results. For further details on approximation methods for semi-separable matrices, and the theoretical background behind them, we refer the reader to (Patrick Dewilde and A.-J. Veen 1998; Melchior, Van Dooren, and Gallivan 2014)

C.2 Causal Low-rank Matrix Approximation

Given a set of self-attention matrices, we tried to find how close an causal low-rank matrix could approximate $M = \text{Softmax}(\mathbf{Q}\mathbf{K}^\top)$. To ensure the state size N , or in this case rank of N , of the approximation \tilde{M} , we composed $\tilde{\mathbf{M}} = \mathbf{L} \circ \mathbf{A}\mathbf{B}^\top$ where $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{D, N}$, \mathbf{L} is a $\mathbb{R}^{D, D}$ lower triangular mask, and $D = 512$.

We used the results from our causal low-rank (LR) experiments to inform much of our experimental design for later

gradient descent-based approximations, which include both SSD classes (with and without \mathbf{D} matrix) and RetNet. We experimented with various different low-rank approximation solvers. We found that gradient descent performed better than alternating gradient descent. Both types of gradient descent were better than alternating least-squares which often times reached less than optimal local minima. Causal low-rank matrix approximation can also be seen as a softer version of the low-rank matrix completion problem, but a semi-definite programming (SDP) approach was not able to outperform standard gradient descent.

Due to our LR approximation requiring gradient descent, we selected the number of steps in relation to the time required to calculate the semi-separable approximation of the same matrix. Given the heuristic approach for converting self-attention matrices to a semi-separable form (Appendix C.1) and its ability to be parallelized, we selected the number of steps for gradient descent based on the time it took to run an entire batch of matrices (32) using gradient descent on causal low-rank versus one matrix using the semi-separable heuristic. After testing with the state sizes $N = 16, 32, 64$, we found that 10,000 steps suitable as it was around a factor of $5\times$ compared to SSM. The 10,000 steps was maintained across all gradient based approximation classes (SSD, SSD without \mathbf{D} , and RetNet). Experiments using the finalized step count showed AdamW provided better results compared to SGD/Adam, and the use of a scheduler provided little gain.

During the experiments, we also found that initialization of the matrices \mathbf{A}, \mathbf{B} played a significant role in the resulting approximation difference. The original \mathbf{A}, \mathbf{B} values were sampled from $[0, 1)$; however, given $\forall M_{ij} \leq 1, i, j \in [D]$ due to the SoftMax operator, \mathbf{A}, \mathbf{B} values was then sampled from $\left[0, \frac{1}{\sqrt{512N}}\right)$ to have the last row of the self-attention matrix be uniform probability. We then proceeded to vary the factor of the range exponentially, testing $\left[0, \frac{1}{\sqrt{512N}}\right) * 2^{\{-2, -1, 0, 1, 2, 4, 8\}}$ where we found $\left[0, \frac{1}{\sqrt{512N}}\right) * 2^4$ provided the best initialization across multiple datasets. A normal distribution with $\mu = 0$ and σ^2 with the above tested values performed worse than the uniform distribution. Initialization experiments were conducted using the AdamW optimizer with a learning rate of 0.001 and the standard 10,000 steps. This and subsequent gradient descent classes use the same initialization for their \mathbf{A}, \mathbf{B} matrices.

For all gradient descent experiments in Table 9, Three learning rates 0.1, 0.01, 0.001 and AdamW were used for each combination of matrix class, state size, and dataset, with the best approximation being documented. The Frobenius matrix norm was used as the loss function.

Listing 2 PyTorch example for generating Causal Low-rank approximation.

```
n_states = 512
state_size = 16 # or 32, 64
num_heads = 32

A = torch.rand((num_heads, n_states, state_size))
B = torch.rand((num_heads, state_size, n_states))

L = torch.tril(torch.ones((n_states, n_states)))

M_approximation = L * (A @ B)
```

C.3 State Space Dual (SSD) Approximation

For the SSD approximation, we utilize the scalar SSM recurrent (1SS) representation introduced in Dao and A. Gu (2024). A key component is the values of a , which we will refer to l from here on out to avoid confusion with matrix \mathbf{A} , that constitute the final matrix mixer \mathbf{M} .

Given the rolling multiplicative property of L and the size of `n_states`, initialization of l was important to prevent the bottom-right values of L quickly reaching 0. We explored the uniform initialization of $[0, 1) + \{-10, -8, -6, -4, -2, 0, 2\}$ where smaller values of l leads to less “decay” within the L matrix. We found sampling l from $[-8, -7)$ resulted in the best performance and use this initialization in the SSD family and RetNet class. As expected, adding the \mathbf{D} component helps reduce the error between the approximation and actual attention matrix Table 9.

Listing 3 PyTorch example for generating SSD (with and without **D** component) approximation.

```
n_states = 512
state_size = 16 # or 32, 64
softplus = torch.nn.Softplus()
num_heads = 32

A = torch.rand((num_heads, n_states, state_size))
B = torch.rand((num_heads, state_size, n_states))
D = torch.rand((num_heads))

l = torch.rand((h, n_states))
L = torch.exp(segsum(softplus(l) * -1))

M_approximation = L * (A @ B)
if apply_D:
    M_approximation = M_approximation + torch.eye(n_states,n_states) * D[:, None, None]
```

C.4 RetNet Matrix Approximation

The Retention mechanism, introduced by Sun et al. (2023), is a key component in RetNet models and can be represented mathematically as $(\mathbf{QK}^T \cdot \mathbf{L})\mathbf{V}$. Here, the matrix \mathbf{L} is defined element-wise by

$$L_{nm} = \begin{cases} \gamma^{n-m}, & n \geq m \\ 0, & n < m \end{cases} \quad (7)$$

where γ is a decay factor. This lower triangular matrix \mathbf{L} captures the temporal dependencies by decaying past values with respect to the current position.

In our approximation, we replace the product \mathbf{QK} with matrices \mathbf{A} and \mathbf{B} . The matrix \mathbf{L} can be efficiently constructed in PyTorch using the following code, which generates a RetNet approximation:

Listing 4 PyTorch example for generating the RetNet matrix approximation.

```
n_states = 512
state_size = 16 # or 32, 64
softplus = torch.nn.Softplus()
num_heads = 32

A = torch.rand((num_heads, n_states, state_size))
B = torch.rand((num_heads, state_size, n_states))

l = torch.rand((num_heads))
L = torch.exp(segsum(softplus(einops.repeat(l, 'h -> h n', n=n_states) * -1))

M_approximation = L * (A @ B)
```

This implementation provides a practical method for simulating the Retention mechanism, crucial for reducing computational complexity in RetNet models.

C.5 Toeplitz Approximation

Our Toeplitz approximation technique calculates the matrix approximation by setting the value of each band of the Toeplitz matrix as the average of the values of the respective band in the attention matrix. Since each band in a Toeplitz matrix is constant along its diagonal, this method ensures that the approximation preserves the structure of the original matrix while maintaining computational efficiency.

To justify this approach, we observe that taking the mean per band minimizes the L2 norm (i.e., the sum of squared differences) between the original attention matrix and the approximated Toeplitz matrix. Specifically, for each band, the optimal value that minimizes the L2 difference between the two matrices is the average of the elements in that band. This is because the mean is the value that minimizes the sum of squared deviations for a set of numbers. As such, using

the mean ensures that the approximation is as close as possible to the original matrix in terms of L2 distance, thereby providing a robust and efficient approximation method.

As before, we assume that the approximation is input-dependent, meaning that each attention matrix has its own unique Toeplitz approximation.

C.6 Segsum Operator

The segsum operator computes the sum of elements across specified segments of a matrix, which, as applied in Appendices C.3 and C.4, corresponds to summing over the columns. This operation is crucial for various matrix manipulations, including the computation of the state-space dual (refer to Equation (2)). Below is the Python implementation of the ‘segsum’ operator using PyTorch.

Listing 5 PyTorch implementation of the Segmented Summation (segsum) operator.

```
def segsum(x):
    """Naive segment sum calculation. exp(segsum(A)) produces a 1-SS matrix,
       which is equivalent to a scalar SSM."""
    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), diagonal=0)
    x_segsum = x_segsum.masked_fill(~mask, -torch.inf)
    return x_segsum
```
