
Achilles’ Heel of Mamba: Essential difficulties of the Mamba architecture demonstrated by synthetic data

Tianyi Chen^{1,2*}, Pengxiao Lin^{1,2*}, Zhiwei Wang^{1,2}, Zhi-Qin John Xu^{1,2,3†}

¹ Institute of Natural Sciences, MOE-LSC, Shanghai Jiao Tong University

² School of Mathematical Sciences, Shanghai Jiao Tong University

³ Shanghai Seres Information Technology Co., Ltd, Shanghai 200040, China

Abstract

State Space Models (SSMs) have emerged as promising alternatives to attention mechanisms, with the Mamba architecture demonstrating impressive performance and linear complexity for processing long sequences. However, the fundamental differences between Mamba and Transformer architectures remain incompletely understood. In this work, we use carefully designed synthetic tasks to reveal Mamba’s inherent limitations. Through experiments, we identify that Mamba’s nonlinear convolution introduces an asymmetry bias that significantly impairs its ability to recognize symmetrical patterns and relationships. Using composite function and inverse sequence matching tasks, we demonstrate that Mamba strongly favors compositional solutions over symmetrical ones and struggles with tasks requiring the matching of reversed sequences. We show these limitations stem not from the SSM module itself but from the nonlinear convolution preceding it, which fuses token information asymmetrically. These insights provide a new understanding of Mamba’s constraints and suggest concrete architectural improvements for future sequence models.

1 Introduction

Large Language Models (LLMs) have achieved remarkable progress and are now widely applied across a broad range of fields (Vaswani et al., 2017; Liu et al., 2018; Devlin et al., 2018; Radford et al., 2019; Touvron et al., 2023; OpenAI, 2023; Brown et al., 2020; Dong et al., 2022; Garg et al., 2022; Trinh et al., 2024; Davies et al., 2021). The performance and inductive biases of such models are largely determined by their underlying architectures. Among these, the Transformer (Vaswani et al., 2017) has become a dominant backbone in LLMs. Its attention mechanism is central to its strong performance (Brown et al., 2020; Olsson et al., 2022; Wang et al., 2024). However, this mechanism also represents one of its major limitations: the computational complexity of attention scales quadratically with sequence length, making Transformers inefficient for long-sequence tasks.

To address this issue, various Transformer variants have been proposed to reduce the computational cost of attention, including sparse attention (Child et al., 2019) and linear attention mechanisms (Katharopoulos et al., 2020). Mamba (Gu and Dao, 2024; Dao and Gu, 2024), which incorporates State Space Model (SSM), has recently garnered significant attention due to its linear complexity with respect to sequence length and its superior performance on long-sequence problems. As a result, SSM offering a natural, computation-efficient alternative similar to linear attention have become a focal point of research.

*Equal contribution.

†Corresponding author: xuzhiqin@sjtu.edu.cn.

To better understand the behavior of large-scale models and to guide meaningful architectural improvements, it is crucial to investigate the underlying causes of the differences between Mamba and Transformer. Given the inherent complexity of natural language tasks, this study leverages simple but meticulously crafted synthetic data.

When examining architectural details more closely, is the difference between Mamba and Transformer limited merely to the replacement of attention with State Space Models (SSMs)? In fact, beyond the use of SSMs, Mamba exhibits several fundamental differences from the Transformer architecture. One of the most critical distinctions lies in Mamba’s use of nonlinear convolution (O’Shea and Nash, 2015). On one hand, the nonlinear convolution enables Mamba to propagate information within a sequence without relying solely on the SSM unlike Transformer, where Transformer depends entirely on attention for intra-sequence communication. On the other hand, this convolution fuses token-level information within the sequence, and the SSM operates on this fused representation to perform matching and extraction.

Notably, the nonlinear convolution in Mamba introduces an intrinsic asymmetry due to the asymmetric structure of convolution kernels. This asymmetry is transferred to the fused token representations and consequently affects how information is extracted. To better understand this limitation, we devised a composite function task (Zhang et al., 2024a,b, 2025) aimed at evaluating how Mamba handles compositional structures.

The composite function task admits two possible solutions: a composite solution and a symmetric solution. We observe that Mamba exhibits a strong bias toward the composite solution while struggling to learn the symmetric one. We find that Mamba struggles to match sequences under order changes—for example, "1234" vs. "4321". To test this limitation, we designed an inverse sequence matching task, where the model must match a sequence with its reversed counterpart. Experimental results confirm that Mamba has difficulty completing this task, whereas Transformer handles it with ease. We further introduced a residual connection, complemented by positional encoding, allowing the SSM to directly handle data without relying on convolution. This modification led to a significant improvement in Mamba’s performance on the inverse sequence matching task, confirming that the core issue lies not within the SSM itself, but rather in the nonlinear convolution outside the SSM. This work therefore provides a new angle of symmetry to understand fundamental mechanisms of Mamba structure.

The main contributions of this work are as follows:

1. We conduct an in-depth analysis of how Mamba solves the composite function task, revealing a fundamental difference in information acquisition between Mamba and Transformer. We show that Mamba tends to rely on convolution to retrieve relevant information.
2. Through systematic experiments and observations, we identify Mamba’s bias toward asymmetric solutions in the composite function task. We demonstrate that this behavior distinguishes Mamba from Transformer and trace the root cause to the asymmetry introduced by convolution.
3. To further examine this asymmetry bias, we design an inverse sequence matching task, in which Mamba exhibits clear difficulties. We show that these difficulties can be effectively addressed through targeted modifications inspired by our earlier findings, leading to substantial performance improvements.
4. By highlighting how nonlinear convolution-based fusion induces an inherent asymmetry in Mamba, we provide insightful guidance for future model improvements and the design of new architectures.

2 Related work

State Space Models (SSM), Mamba, and Their Shortcomings. State space models (SSM) originate from neuromorphic spiking models (Voelker, 2019; Voelker et al., 2019) and have gained prominence through several key developments, such as S4 (Gu et al., 2022), S5 (Smith et al., 2023), the RWKV series (Peng et al., 2023, 2024), RetNet (Sun et al., 2023), and GLA (Yang et al., 2024). Among these, Mamba2 (Gu and Dao, 2024) stands out by delivering performance competitive with Transformers while requiring significantly lower computational resources. However, numerous studies reveal that Mamba has notable limitations. Research including (Ben-Kish et al., 2025; Ye et al.,

2025; Yuan et al., 2025) indicates that Mamba generally underperforms compared to Transformers in tasks involving long-context understanding, prompting the development of alternative models like DeciMamba, LongMamba, and ReMamba. Similarly, (Park et al., 2024; Waleffe et al., 2024) explore Mambas in-context learning abilities and conclude that they fall short of Transformer capabilities. Additionally, (Arora et al., 2023; Jelassi et al., 2024) demonstrate that Mamba struggles with retrieval tasks, such as copying information from the input context. Some studies, such as (Xu et al., 2025), also examine the practical efficiency of Mamba compared to Transformers, noting that despite its theoretical advantages, real-world performance can be lacking. Furthermore, (Ren et al., 2024) introduces a COPY task, which exposes a performance bottleneck in Mamba. Our study takes a structural view of Mamba, employing carefully designed synthetic data to explore its behavior from the perspective of symmetry, offering deeper insights into these shortcomings of Mamba and proposing a solution to address them.

Understanding the Mechanism of Neural Network Models. Our work conducts an in-depth investigation into the internal mechanisms of Mamba. In identifying the functional roles of key modules in Mamba presented in this paper, we adopt the commonly used techniques of perturbation and causal intervention similar to research on interpretability in large language models (Vig et al., 2020; Jeoung and Diesner, 2022; Wang et al., 2022; Conmy et al., 2023; Merullo et al., 2023; Guo et al., 2023; Wang and Weinan, 2024; Amsel et al., 2024; Li et al., 2024; Wang et al., 2024). The construction of synthetic datasets in our study is primarily inspired by the works of Poli et al. (2024) and Zhang et al. (2024a). The design of our composite function tasks is adopted from the approach used in Zhang et al. (2024b,a). Additionally, several insightful theoretical studies on feedforward neural networks have also informed our theoretical analysis. For instance, various works explore the preferences and generalization capabilities of neural networks from perspectives such as regularization, frequency, and dynamics (Xu et al., 2025a,b, 2019; Wang et al., 2024; Jacot et al., 2018, 2020; Arora et al., 2019, 2018; Wu and Su, 2023; Wang and Wu, 2023; Arora et al., 2022; Li et al., 2021; Wu et al., 2018; Zhu et al., 2018; Arora et al., 2019; Ren et al., 2024).

3 Preliminary

3.1 Introduction of Mamba

The overall structure of Mamba (Gu and Dao, 2024; Dao and Gu, 2024) is shown in Fig. 1. The Mamba block can be divided into three parts: the widely known SSM component, the pre-SSM part, and the post-SSM part. Omitting trivial dimension transformations and setting the batch size to 1 to omit the batch dimension, for a given input U to the block, the internal computation process to obtain the output O can be described as follows:

Pre-SSM

$$(\tilde{U}, Z, dt) = \mathbf{Linear}(U), \quad U \in R^{(s,d)}, \tilde{U} \in R^{(s,2d+2h)}, Z \in R^{(s,2d)}, dt \in R^{(s,N_h)}, \quad (1)$$

$$(B, C, X) = \sigma(\mathbf{Conv1d}(\tilde{U})), \quad B \in R^{(s,h)}, C \in R^{(s,h)}, X \in R^{(s,2d)}, \quad (2)$$

$$\tilde{X} = X \circ dt, X \in R^{(N_h,s,2d/N_h)}, \quad (3)$$

SSM

$$Mask = \mathbf{F}(A_0, dt), \quad A_0 \in R^{N_h}, Mask \in R^{(N_h,s,s)}, \quad (4)$$

$$I = \mathbf{Repeat}(CB^\top, N_h), \quad I \in R^{(N_h,s,s)}, \quad (5)$$

$$S = Mask \circ I, \quad S \in R^{(N_h,s,s)}, \quad (6)$$

$$Y = S\tilde{X} + X, \quad Y \in R^{(N_h,s,2d/N_h)}, \quad (7)$$

Post-SSM

$$Y_{Norm} = \mathbf{RMS}(Y \circ (\sigma(Z))), \quad Y_{Norm} \in R^{(s,2d)}, \quad (8)$$

$$O = \mathbf{Linear}(Y_{Norm}), \quad O \in R^{(s,2d)}, \quad (9)$$

where s is the sequence length, d is the model’s hiddenstate dimension, h is the SSM hidden dimension, S is the SSM matrix, N_h is the number of SSM heads, **Linear** denotes a linear transformation,

Conv1d denotes a one-dimensional convolution, σ denotes a nonlinear activation function, **F** denotes the function that generates the *Mask*, **Repeat** denotes a dimension replication operation, \circ denotes pointwise multiplication, and **RMS** denotes RMS normalization.

We define the composition of σ and **Conv1d** as **nonlinear convolution** and **attention score** from token j to token i is given by the (i, j) entry of the SSM matrix S throughout the paper. For more comprehensive computational details, please refer to the appendix.

3.2 Difference between Mamba and Transformer

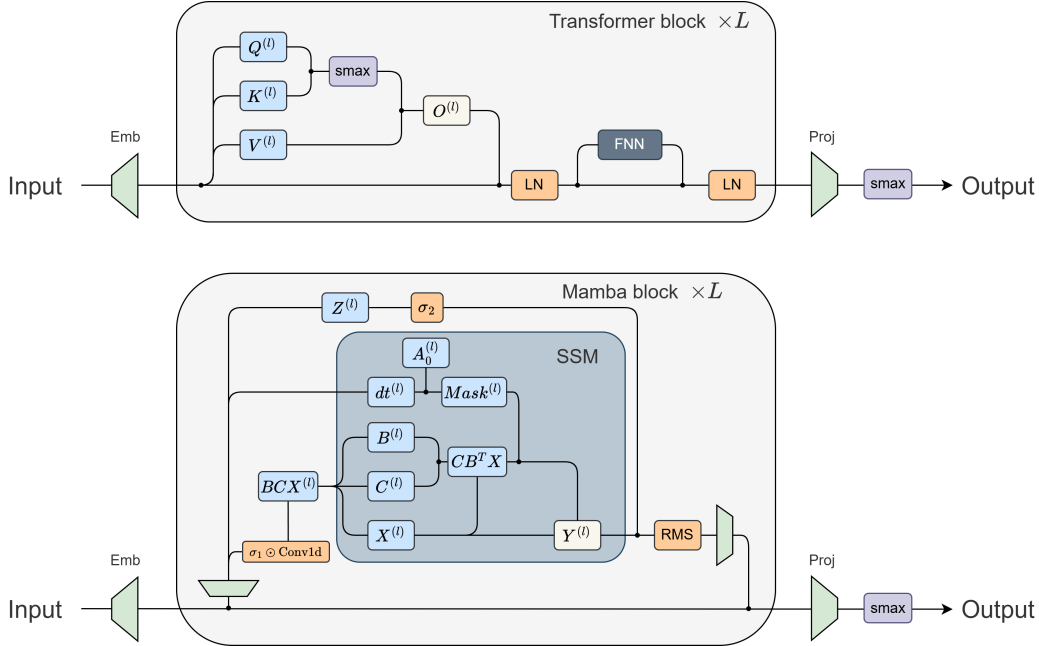


Figure 1: Overview of Mamba and Transformer Blocks. The green trapezoids represent linear mappings. "smax" denotes the softmax function, "FNN" stands for feed-forward neural network, and "LN" represents layer normalization. The meanings of variables specific to the Mamba block are explained in the main text.

By incorporating the SSM module, Mamba circumvents the quadratic complexity of attention in Transformers. While SSM significantly improves computational efficiency, the key difference from attention appears limited to the softmax function and the learnable mask. However, notable distinctions also exist beyond the SSM itself.

Prior to the attention-like SSM, Mamba applies a nonlinear convolution that fuses information, limiting the SSM to operate on already mixed representations. Mamba contains a z-gate following SSM, lacks a Feed-Forward Network (FFN), and its preceding nonlinear layer lacks hidden layers.

As we will demonstrate, nonlinear convolution is a defining feature of Mamba that underlies its fundamental divergence from the Transformer.

3.3 Composite function task

A detailed illustration of the composite function task (Zhang et al., 2024a,b, 2025) is shown in Fig. 2. The core idea of the composite function task can be understood with a simple real world analogy. Imagine a row of people indexed by numbers. Let's define two functions: $n(x)$: Find the person n positions to the right of person x and $m(x)$: Find the person m positions to the left of person x . The composite function $m(n(x))$ means "If you start with person x , go n people to their right, and then m people to their left, who do you end up with?"

In the composite function task, each sequence contains two anchors and one key, the remaining elements in the sequence are randomly sampled from the same range as the key. Each anchor is

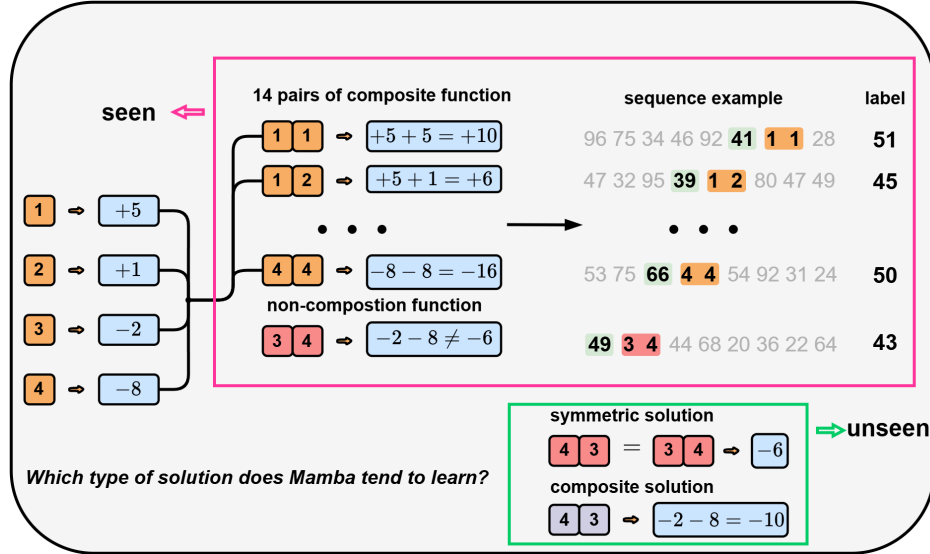


Figure 2: Illustration of the Composite Function Task. Anchors 1, 2, 3, and 4 (depicted in orange) each represent distinct functions. Among the 16 anchor pairs formed, 14 correspond to composite functions derived directly from the sequential application of the individual anchor functions. The pair "34", highlighted in red, is defined as a different function rather than a direct composition. The pair "43" is intentionally excluded from the training set. The input to each anchor pair function is referred to as the "key" (indicated in green). Label indicates the output of an anchor pair applied to a key.

assigned a unique function, and a pair of anchors defines a composite function, with the key (token before anchor pairs) serving as the input. We use the token 1 through 4 as anchors, yielding a total of 16 possible anchor pairs. Among these, 14 composite functions are defined by directly composing the functions associated with the individual anchors. For the anchor pair "34", however, we assign a function that deviates from the standard composition of its components.

The objective is to examine how Mamba handles the unseen anchor pair "43". Two plausible solutions exist: one is the **symmetric solution**, which infers the result of "43" by symmetry from the result of "34"; the other is the **composite solution**, which computes the result by directly composing the functions associated with "4" and "3".

In all experiments, the loss is calculated exclusively based on the final token of the sequence and its corresponding label. For detailed data and training settings, please refer to the appendix.

3.4 Inverse sequence matching task

To investigate whether Mamba can match token information fused by the nonlinear convolution, we design an inverse sequence matching task. The detailed structure of the task is illustrated in Fig. 3. The inverse sequence matching task corresponds to real world situations where people search for the symmetrical counterpart of an object or text. It is similar to asking children to find the symmetrical counterpart of a toy and fit them together. This task can test a model's ability to perform symmetry matching. Although it is simple and fundamental, it reflects a core capability of the model.

Each sample is built from a generating set of three distinct numbers (e.g., $\{28, 92, 37\}$). From its six possible permutations, five are randomly chosen to form key sequences, which are concatenated and separated by a token not in the generating set.

Next, one of the five key sequences is randomly selected, reversed, and appended to the end as the query sequence, separated by non-generating-set tokens so that the query and key sequences would not fuse by convolution. The number of tokens is determined by Mamba's pure convolutional receptive field. For example, in the illustration, six tokens correspond to the pure convolutional receptive field (refers exclusively to how many tokens can be accessed through convolution alone,

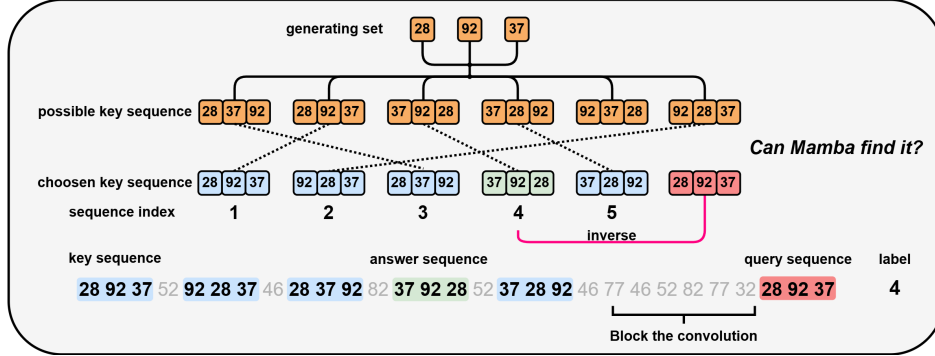


Figure 3: Illustration of the Inverse Sequence Matching Task. The orange elements denote the generation set, which consists of three distinct numbers randomly selected from the interval $[20, 100]$, as well as all possible permutations thereof. Blue and green indicate selected key sequences from the permutation space, separated by random numbers that do not belong to the generation set. One green key sequence is chosen as the answer sequence. The query sequence (shown in red) is obtained by reversing the answer sequence. The corresponding label identifies the position of the answer sequence. To prevent Mamba from leveraging its nonlinear convolution mechanism to infer answers, we prepend the query sequence with random numbers (outside the generation set) matching the length of Mamba’s pure convolutional receptive field.

without invoking the SSM) of a two-layer Mamba ($2 \times 3 = 6$). Finally, the label is the position index of the answer sequence, i.e., the unreversed version of the query sequence.

4 Mamba biases composite solution

In this section, we will empirically show that Mamba is struggling with the symmetric property in composition but biases composite solutions, and analyze the role of nonlinear convolution.

Initialization scale is shown to be critical for a Transformer to learn the composite or symmetric solution (Zhang et al., 2024b, 2025). Therefore, we scan different initialization scales as follows. A parameter $W \in \mathbb{R}^{d_1 \times d_2}$ is initialized as a Gaussian distribution $\mathcal{N}(0, (1/d_1^\gamma)^2)$, where γ is called initialization rate (Luo et al., 2021).

We scan a Mamba with different layers and with different γ . As shown in Fig. 4, where the accuracy is computed by regarding the label as the composite solution in Fig. 4a or the symmetric solution in Fig. 4b, for small γ (large initialization), the network fails to capture either composite or symmetric solution. The value at each position in the figure is computed as the average of three independent random runs. This is consistent with previous study in Transformer (Zhang et al., 2024b, 2025). For large γ (small initialization), the network biases the composite function across almost all cases.

4.1 SSM does not function in composite tasks

It requires information from anchors and key to finish a composition task. Mamba has two potential options: utilizing the convolution or the SSM module. We found that in standard Mamba (Dao and Gu, 2024), convolution plays a critical role while SSM module does not function.

Information blocking. The information flow analysis is a useful tool to visualize the information exchange on token level (Wang et al., 2024). Since the SSM module has high similarity with the attention, we treat every element in SSM matrix as the “attention score” in information flow in Fig. 5. The result suggests that, within the SSM, tokens at later positions have little attention to the key and anchor information. In the case of the Transformer, failure to retrieve the required information via attention renders it incapable of solving the composite function task. To further verify that Mamba does not utilize the SSM for information propagation, we applied a causal intervention approach (Feng and Steinhardt, 2023; Meng et al., 2022; Vig et al., 2020; Wang et al., 2024), manually blocking all information flow from the key and both anchors to the downstream tokens. The (i, j) element of S represents how much token i attends to token j . If this value is set to 0, it implies that token

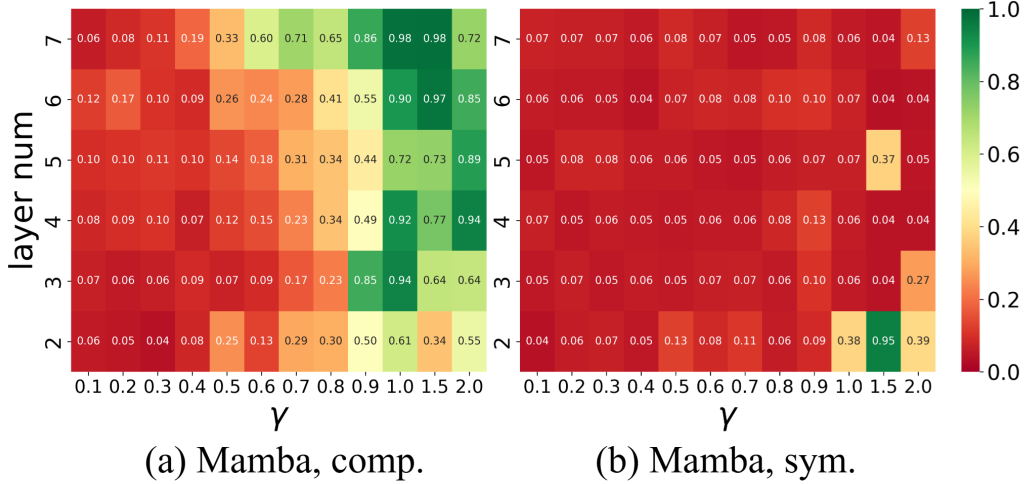


Figure 4: Phase diagram of Mamba on the composite function task. Accuracy (color) for composite function task under different initialization rates (abscissa) and depths (ordinate). The groundtruth for (a) is composite solution and for (b) is symmetric solution. Detailed model configurations and training settings are provided in the appendix.

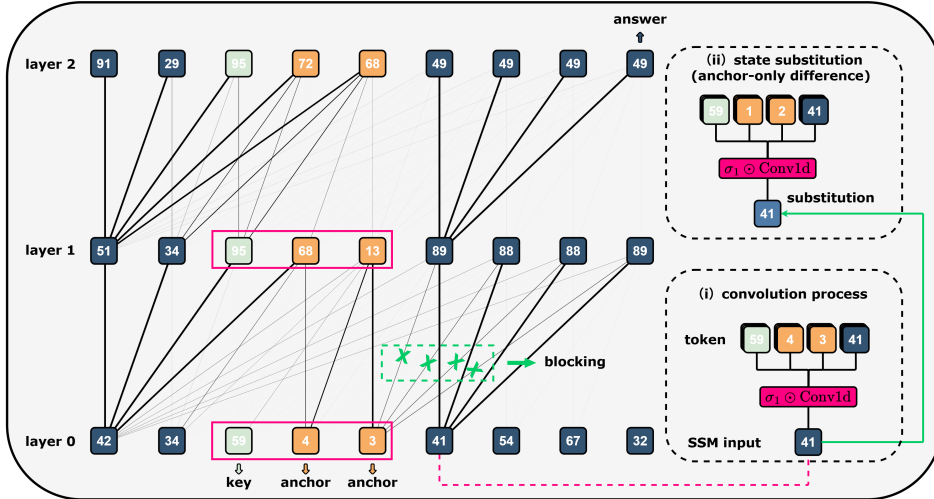


Figure 5: SSM information flow in Mamba for the composite function task. Left: SSM information flow; green crosses indicate pruned connections. Right: (i) SSM input computation; (ii) state replacement after convolution. Flow is computed from $S = Mask \circ C^T B$, with line thickness indicating flow magnitude. Attention score from token j to token i is given by the (i, j) entry of S . The numbers are the outputs of each layer through the model's final linear layer and then take the argmax of the resulting logits to obtain the corresponding digit.

i cannot receive information from token j through the SSM. Our blocking mechanism is implemented by zeroing out the specific entries in S corresponding to the connections we wish to block. As shown in the Fig. 6, for various anchor pairs, the output after cutting these connections remains nearly identical to the original output, indicating that the SSM plays little role in transmitting this information.

Information substitution. To further verify that Mamba relies on convolution for information transmission, we conducted an information substitution experiment: if Mamba encodes all necessary information through convolution, then transferring the resulting state to another sequence should enable it to produce the same output as the original. For all sequences with anchor pairs other

than ‘43’, as illustrated in Fig. 5, we replaced the post-convolution hidden states of the downstream tokens with those from a “43” sequence, where all other elements are identical except for the anchor pair. As shown in Fig. 6, we found that the outputs of nearly all anchor pairs collapse to the output corresponding to the “43” anchor pair.

Taken together, these two experiments clearly demonstrate that Mamba solves the composite function task primarily by leveraging convolution to extract the necessary information. This lays the groundwork for Mamba’s inability to reach symmetric solution.

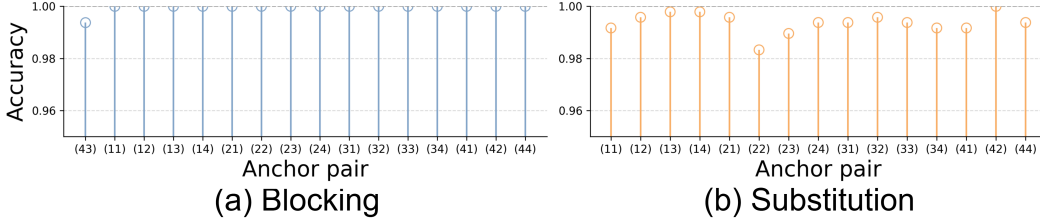


Figure 6: Accuracy of each anchor pair under blocking and substitution experiments compared to the original output. Figures (a) and (b) correspond to the blocking and substitution experiments, respectively. The x-axis denotes the anchor pairs, and the y-axis represents accuracy. Each anchor pair is evaluated using 480 randomly generated sequences.

4.2 Nonlinear convolution introduces asymmetry

In this section we will examine that the asymmetry weight in convolution introduces asymmetry between symmetric anchor pairs.

We design a task with fully symmetric setting as follows. The symmetric anchor pairs have the same function but they do not composite by elementary functions; only function “43” is unseen during the training. This leaves the symmetric solution as the only possible one. As shown in the Fig. 7a, we found that standard Mamba can achieve 100% in the test data of “34” function and fails to learn “43” via symmetric property. Additional experimental results are provided in the appendix.

To better isolate the effect of convolution’s asymmetry, we remove this asymmetry by setting all convolutional kernel weights to 1. As shown in the Fig. 7b, once the asymmetry of convolution is eliminated, Mamba biases learning the symmetric solution to fit the composite task.

This indicates that Mamba’s preference for asymmetry stems from the inherent asymmetry of its nonlinear convolution. In Mamba, sequence information is fused through a nonlinear convolution operation, which serves as the input to the SSM module. Consider two sequences, (v_1, v_2, v_3, v_4) and its reversed counterpart. For the convolution outputs of the above sequences, we define the final token of the result as follows:

$$\text{original sequence: } f = c_1 \circ v_1 + c_2 \circ v_2 + c_3 \circ v_3 + c_4 \circ v_4 + \beta,$$

$$\text{symmetric sequence: } g = c_1 \circ v_4 + c_2 \circ v_3 + c_3 \circ v_2 + c_4 \circ v_1 + \beta.$$

If the convolution parameters c_i are not identical (β is the bias), then f and g will generally differ. This means that even for token sequences that are symmetric in content, their representations after convolution in Mamba can be significantly different. Such a discrepancy illustrates Mamba’s inherent asymmetry, as it fails to preserve the equivalence of symmetric inputs.

The root cause of this asymmetry lies in the non-uniformity of the convolution weights. Since the convolution parameters in Mamba are initialized randomly and trained independently, the c_i values typically remain distinct throughout training. As a result, the convolution operation induces a persistent asymmetry, where different token orders lead to different outputs. We examined the cosine similarity between the individual parameters of the convolution kernel at both the beginning and end of training, and found that they were largely orthogonal to one another, indicating a strong and persistent inconsistency throughout the training process. The detailed results can be found in the appendix.

For a Mamba network with initialization rate $\gamma = 0.5$, it cannot learn composite function task by either composite function or symmetric function. We found that if positional encoding is explicitly

included, as shown in the Fig. 7c, such Mamba network learns composite task by symmetric function. Therefore, positional encoding is also critical for learning symmetric solution.

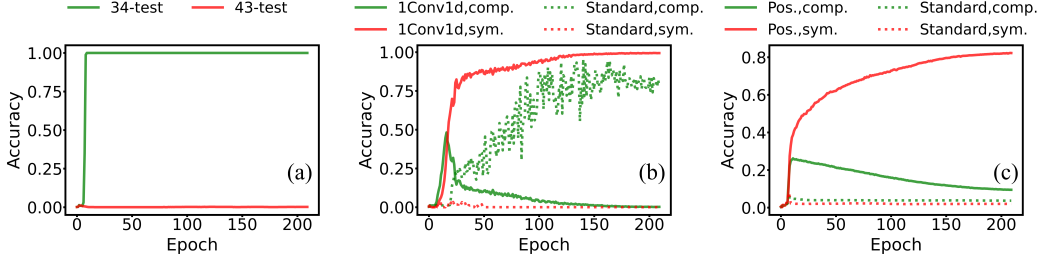


Figure 7: Test Accuracy across. (a): Fully symmetric setting, i.e., target functions are symmetric but do not composite by elementary functions, standard Mamba with $\gamma = 0.5$. (b): Accuracy computed based on groundtruth of composite solution (comp.) and symmetric solution (sym.) for Mamba with all-one convolution (1Conv1d) and standard Mamba (Standard structure) with initialization rate $\gamma = 1$. (c): Similar legend as (b) for Mamba with positional encoding (Pos.) and the standard one but with initialization rate $\gamma = 0.5$. Details of the data and training setup can be found in the appendix.

4.3 Transformer with convolution biases asymmetric solution

The empirical analysis in Mamba reveals that the convolution structure is a key component to introduce asymmetry. Imagine if convolution is introduced into the Transformer, does it increase anchor asymmetry and push the model toward learning asymmetric composite solution? We insert a convolution after the input to the Transformer and before the attention module (applying it to Q , K , and V). This is analogous to how Mamba applies nonlinear convolution before the SSM module. For detailed configurations, please refer to the appendix. Previous study (Zhang et al., 2024b) has shown that with different γ , Transformer can learn symmetric solution or composite solution in different regimes. As shown in Fig. 8, with convolution component, Transformer either can not generalize for small γ or fit data by symmetric solution for relative large γ . In this case, the Transformer exhibits a similar preference for composite function solutions as Mamba and struggles to learn symmetric solutions.

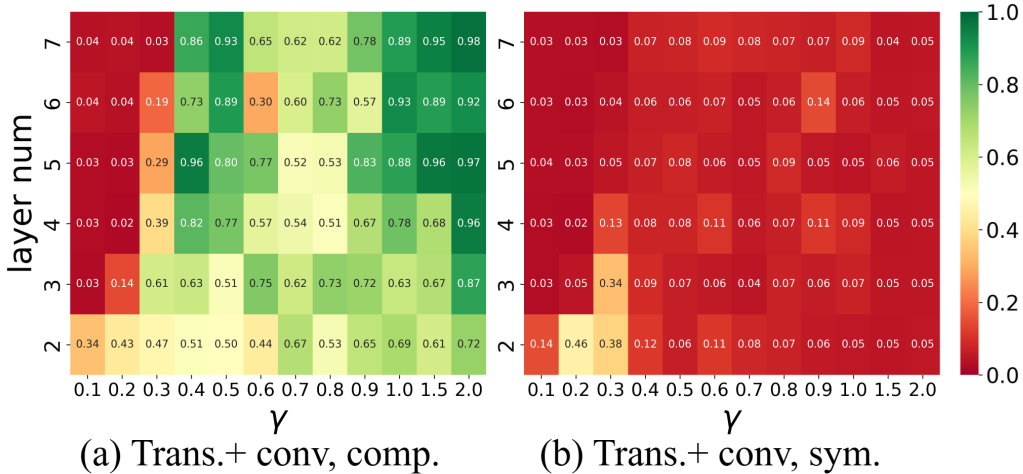


Figure 8: Phase Diagram of the Composite Function Task under Different Settings. Accuracy (color) for composite function task under different initialization rates (abscissa) and depths (ordinate). (a) and (b) show the composite and symmetric solution accuracy of the Transformer after adding non-linear convolution.

5 Mamba’s challenges in the inverse sequence matching task

This raises a critical question: Does the change in token order between sequences like 1234 and 4321 make it inherently difficult for Mamba to attend across reversed patterns and discover the correct reverse sequence?

As shown in Fig. 9a, for the standard Mamba network, the training accuracy can easily achieve 100%, however, the test accuracy for the case with tokens seen during training or the case with tokens unseen during the training (“OOD”), the network predicts the outcome with a random-guess level. Additional experimental results can be found in the appendix. For Transformer, or a Mamba network that adds a residual connection from the input of convolution to SSM and the position embedding, as shown in Fig. 9b and 9c, the network can accurately predicts test cases. In addition, for the “OOD” case, the accuracies of such two cases are also significantly larger than random guess. We also conducted additional architectural experiments aimed at mitigating Mamba’s asymmetry bias. For details, please refer to the appendix.

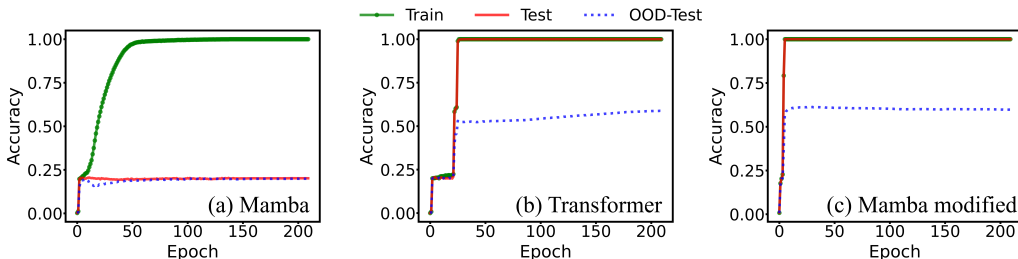


Figure 9: Accuracy on the inverse sequence matching task across different model architectures. (a), (b), and (c) show the accuracy of Mamba, Transformer, and the modified Mamba with residual connection, respectively, on the inverse sequence matching task under the same settings. “OOD” refers to a set drawn from a distribution outside the training and standard test ranges.

6 Conclusion

In this work, we leveraged synthetic data to experimentally and systematically analyze the intrinsic properties of the Mamba architecture. Our findings reveal fundamental differences between Mamba and Transformer models, particularly in handling symmetrical patterns and relationships. We identified that Mamba’s inherent asymmetry bias stems from its nonlinear convolution mechanism, which fuses token information asymmetrically before passing it to the SSM module. This architectural constraint limits Mamba’s ability to recognize symmetrical solutions in composite function tasks and to match reversed sequences. These insights provide valuable guidance for designing future sequence models that combine the computational efficiency of SSMs with the flexible pattern recognition capabilities of Transformers.

7 Limitations

To enable a precise investigation and clear illustration of Mamba’s internal mechanisms and inductive biases, this work employs synthetic data rather than real-world datasets. While the synthetic tasks are designed to capture key characteristics of the Mamba architecture, they may not generalize to the full diversity of real-world data. Moreover, due to the use of synthetic data, the Mamba models used in our experiments are relatively small. Whether the observed biases persist in larger-scale Mamba architectures remains an open question requiring further investigation. Additionally, many existing large-scale models incorporating Mamba do so in conjunction with Transformer components. In such hybrid architectures, it remains to be studied whether Mamba’s inductive bias still dominates or is diminished.

Acknowledgements

This work is supported by the National Key R&D Program of China Grant No. 2022YFA1008200, the National Natural Science Foundation of China Grant No. 92270001, 12371511, 12422119, 2025 Key Technology R&D Program “New Generation Information Technology” Project of Shanghai Municipal Science and Technology Commission, and the HPC of School of Mathematical Sciences and the Student Innovation Center, and the Siyuan-1 cluster supported by the Center for High Performance Computing at Shanghai Jiao Tong University, and Key Laboratory of Marine Intelligent Equipment and System (Ministry of Education, P.R. China), and SJTU Kunpeng & Ascend Center of Excellence.

References

- A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, I. Polosukhin, Attention is all you need, *Advances in neural information processing systems* 30 (2017).
- P. J. Liu, M. Saleh, E. Pot, B. Goodrich, R. Sepassi, L. Kaiser, N. Shazeer, Generating wikipedia by summarizing long sequences, *arXiv preprint arXiv:1801.10198* (2018).
- J. Devlin, M.-W. Chang, K. Lee, K. Toutanova, Bert: Pre-training of deep bidirectional transformers for language understanding, *arXiv preprint arXiv:1810.04805* (2018).
- A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, I. Sutskever, Language models are unsupervised multitask learners (2019).
- H. Touvron, T. Lavril, G. Izacard, X. Martinet, M.-A. Lachaux, T. Lacroix, B. Rozière, N. Goyal, E. Hambro, F. Azhar, et al., Llama: Open and efficient foundation language models, *arXiv preprint arXiv:2302.13971* (2023).
- OpenAI, Gpt-4 technical report, 2023. *arXiv:2303.08774*.
- T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, et al., Language models are few-shot learners, *Advances in neural information processing systems* 33 (2020) 1877–1901.
- Q. Dong, L. Li, D. Dai, C. Zheng, Z. Wu, B. Chang, X. Sun, J. Xu, Z. Sui, A survey on in-context learning, *arXiv preprint arXiv:2301.00234* (2022).
- S. Garg, D. Tsipras, P. S. Liang, G. Valiant, What can transformers learn in-context? a case study of simple function classes, *Advances in Neural Information Processing Systems* 35 (2022) 30583–30598.
- T. H. Trinh, Y. Wu, Q. V. Le, H. He, T. Luong, Solving olympiad geometry without human demonstrations, *Nature* 625 (2024) 476–482.
- A. Davies, P. Velikovi, L. Buesing, S. Blackwell, D. Zheng, N. Tomaev, R. Tanburn, P. Battaglia, C. Blundell, A. Juhász, et al., Advancing mathematics by guiding human intuition with ai, *Nature* 600 (2021) 70–74.
- C. Olsson, N. Elhage, N. Nanda, N. Joseph, N. DasSarma, T. Henighan, B. Mann, A. Askell, Y. Bai, A. Chen, et al., In-context learning and induction heads, *arXiv preprint arXiv:2209.11895* (2022).
- Z. Wang, Y. Wang, Z. Zhang, Z. Zhou, H. Jin, T. Hu, J. Sun, Z. Li, Y. Zhang, Z.-Q. J. Xu, The buffer mechanism for multi-step information reasoning in language models, *arXiv preprint arXiv:2405.15302* (2024).
- R. Child, S. Gray, A. Radford, I. Sutskever, Generating long sequences with sparse transformers, *arXiv preprint arXiv:1904.10509* (2019).
- A. Katharopoulos, A. Vyas, N. Pappas, F. Fleuret, Transformers are rnns: Fast autoregressive transformers with linear attention, in: *International conference on machine learning*, PMLR, 2020, pp. 5156–5165.

- A. Gu, T. Dao, Mamba: Linear-time sequence modeling with selective state spaces, The Conference on Language Modeling (2024).
- T. Dao, A. Gu, Transformers are ssms: Generalized models and efficient algorithms through structured state space duality, arXiv preprint arXiv:2405.21060 (2024).
- K. O’shea, R. Nash, An introduction to convolutional neural networks, arXiv preprint arXiv:1511.08458 (2015).
- Z. Zhang, Z. Wang, J. Yao, Z. Zhou, X. Li, Z.-Q. J. Xu, et al., Anchor function: a type of benchmark functions for studying language models, arXiv preprint arXiv:2401.08309 (2024a).
- Z. Zhang, P. Lin, Z. Wang, Y. Zhang, Z.-Q. J. Xu, Initialization is critical to whether transformers fit composite functions by inference or memorizing, arXiv preprint arXiv:2405.05409 (2024b).
- Z. Zhang, P. Lin, Z. Wang, Y. Zhang, Z.-Q. J. Xu, Complexity control facilitates reasoning-based compositional generalization in transformers, arXiv preprint arXiv:2501.08537 (2025).
- A. R. Voelker, Dynamical systems in spiking neuromorphic hardware (2019).
- A. Voelker, I. Kaji, C. Eliasmith, Legendre memory units: Continuous-time representation in recurrent neural networks, Advances in neural information processing systems 32 (2019).
- A. Gu, K. Goel, C. Re, Efficiently modeling long sequences with structured state spaces, in: International Conference on Learning Representations, 2022. URL: <https://openreview.net/forum?id=uYLFoz1v1AC>.
- J. T. Smith, A. Warrington, S. Linderman, Simplified state space layers for sequence modeling, in: The Eleventh International Conference on Learning Representations, 2023. URL: <https://openreview.net/forum?id=Ai8Hw3AXqks>.
- B. Peng, E. Alcaide, Q. Anthony, A. Albalak, S. Arcadinho, S. Biderman, H. Cao, X. Cheng, M. Chung, L. Derczynski, X. Du, M. Grella, K. Gv, X. He, H. Hou, P. Kazienko, J. Kocon, J. Kong, B. Koptyra, H. Lau, J. Lin, K. S. I. Mantri, F. Mom, A. Saito, G. Song, X. Tang, J. Wind, S. Woźniak, Z. Zhang, Q. Zhou, J. Zhu, R.-J. Zhu, RWKV: Reinventing RNNs for the transformer era, in: H. Bouamor, J. Pino, K. Bali (Eds.), Findings of the Association for Computational Linguistics: EMNLP 2023, Association for Computational Linguistics, Singapore, 2023, pp. 14048–14077. URL: <https://aclanthology.org/2023.findings-emnlp.936/>. doi:10.18653/v1/2023.findings-emnlp.936.
- B. Peng, D. Goldstein, Q. G. Anthony, A. Albalak, E. Alcaide, S. Biderman, E. Cheah, T. Ferdinan, K. K. GV, H. Hou, S. Krishna, R. M. Jr., N. Muennighoff, F. Obeid, A. Saito, G. Song, H. Tu, R. Zhang, B. Zhao, Q. Zhao, J. Zhu, R.-J. Zhu, Eagle and finch: RWKV with matrix-valued states and dynamic recurrence, in: First Conference on Language Modeling, 2024. URL: <https://openreview.net/forum?id=soz1SEiPeq>.
- Y. Sun, L. Dong, S. Huang, S. Ma, Y. Xia, J. Xue, J. Wang, F. Wei, Retentive network: A successor to transformer for large language models, 2023. URL: <https://arxiv.org/abs/2307.08621>. arXiv:2307.08621.
- S. Yang, B. Wang, Y. Shen, R. Panda, Y. Kim, Gated linear attention transformers with hardware-efficient training, in: R. Salakhutdinov, Z. Kolter, K. Heller, A. Weller, N. Oliver, J. Scarlett, F. Berkenkamp (Eds.), Proceedings of the 41st International Conference on Machine Learning, volume 235 of *Proceedings of Machine Learning Research*, PMLR, 2024, pp. 56501–56523. URL: <https://proceedings.mlr.press/v235/yang24ab.html>.
- A. Ben-Kish, I. Zimmerman, S. Abu-Hussein, N. Cohen, A. Globerson, L. Wolf, R. Giryes, Decimamba: Exploring the length extrapolation potential of mamba, in: The Thirteenth International Conference on Learning Representations, 2025. URL: <https://openreview.net/forum?id=iWS15Zyjjw>.
- Z. Ye, K. Xia, Y. Fu, X. Dong, J. Hong, X. Yuan, S. Diao, J. Kautz, P. Molchanov, Y. C. Lin, Longmamba: Enhancing mamba’s long-context capabilities via training-free receptive field enlargement, in: The Thirteenth International Conference on Learning Representations, 2025. URL: <https://openreview.net/forum?id=fMbLszV01H>.

- D. Yuan, J. Liu, B. Li, H. Zhang, J. Wang, X. Cai, D. Zhao, Remamba: Equip mamba with effective long-sequence modeling, 2025. URL: <https://arxiv.org/abs/2408.15496>. arXiv:2408.15496.
- J. Park, J. Park, Z. Xiong, N. Lee, J. Cho, S. Oymak, K. Lee, D. Papailiopoulos, Can mamba learn how to learn? a comparative study on in-context learning tasks, in: ICLR 2024 Workshop on Mathematical and Empirical Understanding of Foundation Models, 2024. URL: <https://openreview.net/forum?id=xvr0Hctddy>.
- R. Waleffe, W. Byeon, D. Riach, B. Norick, V. Korthikanti, T. Dao, A. Gu, A. Hatamizadeh, S. Singh, D. Narayanan, G. Kulshreshtha, V. Singh, J. Casper, J. Kautz, M. Shoeybi, B. Catanzaro, An empirical study of mamba-based language models, 2024. URL: <https://arxiv.org/abs/2406.07887>. arXiv:2406.07887.
- S. Arora, S. Eyuboglu, A. Timalsina, I. Johnson, M. Poli, J. Zou, A. Rudra, C. Ré, Zoology: Measuring and improving recall in efficient language models, 2023. URL: <https://arxiv.org/abs/2312.04927>. arXiv:2312.04927.
- S. Jelassi, D. Brandfonbrener, S. M. Kakade, E. Malach, Repeat after me: Transformers are better than state space models at copying, 2024. URL: <https://arxiv.org/abs/2402.01032>. arXiv:2402.01032.
- Z. Xu, J. Yan, A. Gupta, V. Srikumar, State space models are strong text rerankers, 2025. URL: <https://arxiv.org/abs/2412.14354>. arXiv:2412.14354.
- R. Ren, Z. Li, Y. Liu, Can mamba always enjoy the "free lunch"?, 2024. URL: <https://arxiv.org/abs/2410.03810>. arXiv:2410.03810.
- J. Vig, S. Gehrmann, Y. Belinkov, S. Qian, D. Nevo, Y. Singer, S. Shieber, Investigating gender bias in language models using causal mediation analysis, *Advances in neural information processing systems* 33 (2020) 12388–12401.
- S. Jeoung, J. Diesner, What changed? investigating debiasing methods using causal mediation analysis, arXiv preprint arXiv:2206.00701 (2022).
- K. Wang, A. Variengien, A. Conmy, B. Shlegeris, J. Steinhardt, Interpretability in the wild: a circuit for indirect object identification in gpt-2 small, arXiv preprint arXiv:2211.00593 (2022).
- A. Conmy, A. Mavor-Parker, A. Lynch, S. Heimersheim, A. Garriga-Alonso, Towards automated circuit discovery for mechanistic interpretability, *Advances in Neural Information Processing Systems* 36 (2023) 16318–16352.
- J. Merullo, C. Eickhoff, E. Pavlick, Circuit component reuse across tasks in transformer language models, arXiv preprint arXiv:2310.08744 (2023).
- T. Guo, W. Hu, S. Mei, H. Wang, C. Xiong, S. Savarese, Y. Bai, How do transformers learn in-context beyond simple functions? a case study on learning with representations, arXiv preprint arXiv:2310.10616 (2023).
- M. Wang, E. Weinan, Understanding the expressive power and mechanisms of transformer for sequence modeling, arXiv preprint arXiv:2402.00522 (2024).
- N. Amsel, G. Yehudai, J. Bruna, On the benefits of rank in attention layers, arXiv preprint arXiv:2407.16153 (2024).
- Z. Li, H. Liu, D. Zhou, T. Ma, Chain of thought empowers transformers to solve inherently serial problems, arXiv preprint arXiv:2402.12875 (2024).
- X. Wang, A. Amayuelas, K. Zhang, L. Pan, W. Chen, W. Y. Wang, Understanding the reasoning ability of language models from the perspective of reasoning paths aggregation, arXiv preprint arXiv:2402.03268 (2024).
- M. Poli, A. W. Thomas, E. Nguyen, P. Ponnusamy, B. Deiseroth, K. Kersting, T. Suzuki, B. Hie, S. Ermon, C. Ré, et al., Mechanistic design and scaling of hybrid architectures, arXiv preprint arXiv:2403.17844 (2024).

- Z.-Q. J. Xu, Y. Zhang, T. Luo, Overview frequency principle/spectral bias in deep learning, *Communications on Applied Mathematics and Computation* 7 (2025a) 827–864.
- Z.-Q. J. Xu, Y. Zhang, Z. Zhou, An overview of condensation phenomenon in deep learning, *arXiv preprint arXiv:2504.09484* (2025b).
- Z.-Q. J. Xu, Y. Zhang, T. Luo, Y. Xiao, Z. Ma, Frequency principle: Fourier analysis sheds light on deep neural networks, *arXiv preprint arXiv:1901.06523* (2019).
- M. Wang, H. He, J. Wang, Z. Wang, G. Huang, F. Xiong, Z. Li, L. Wu, et al., Improving generalization and convergence by enhancing implicit regularization, *arXiv preprint arXiv:2405.20763* (2024).
- A. Jacot, F. Gabriel, C. Hongler, Neural tangent kernel: Convergence and generalization in neural networks, *Advances in neural information processing systems* 31 (2018).
- A. Jacot, B. Simsek, F. Spadaro, C. Hongler, F. Gabriel, Implicit regularization of random feature models, in: *International Conference on Machine Learning*, PMLR, 2020, pp. 4631–4640.
- S. Arora, S. Du, W. Hu, Z. Li, R. Wang, Fine-grained analysis of optimization and generalization for overparameterized two-layer neural networks, in: *International Conference on Machine Learning*, PMLR, 2019, pp. 322–332.
- S. Arora, N. Cohen, N. Golowich, W. Hu, A convergence analysis of gradient descent for deep linear neural networks, *arXiv preprint arXiv:1810.02281* (2018).
- L. Wu, W. J. Su, The implicit regularization of dynamical stability in stochastic gradient descent, in: *International Conference on Machine Learning*, PMLR, 2023, pp. 37656–37684.
- M. Wang, L. Wu, A theoretical analysis of noise geometry in stochastic gradient descent, *arXiv preprint arXiv:2310.00692* (2023).
- S. Arora, Z. Li, A. Panigrahi, Understanding gradient descent on the edge of stability in deep learning, in: *International Conference on Machine Learning*, PMLR, 2022, pp. 948–1024.
- Z. Li, S. Malladi, S. Arora, On the validity of modeling sgd with stochastic differential equations (sdes), *Advances in Neural Information Processing Systems* 34 (2021) 12712–12725.
- L. Wu, C. Ma, et al., How sgd selects the global minima in over-parameterized learning: A dynamical stability perspective, *Advances in Neural Information Processing Systems* 31 (2018).
- Z. Zhu, J. Wu, B. Yu, L. Wu, J. Ma, The anisotropic noise in stochastic gradient descent: Its behavior of escaping from sharp minima and regularization effects, *arXiv preprint arXiv:1803.00195* (2018).
- S. Arora, S. S. Du, W. Hu, Z. Li, R. R. Salakhutdinov, R. Wang, On exact computation with an infinitely wide neural net, *Advances in neural information processing systems* 32 (2019).
- Y. Ren, C. Ma, L. Ying, Understanding the generalization benefits of late learning rate decay, in: *International Conference on Artificial Intelligence and Statistics*, PMLR, 2024, pp. 4465–4473.
- T. Luo, Z.-Q. J. Xu, Z. Ma, Y. Zhang, Phase diagram for two-layer relu neural networks at infinite-width limit, *The Journal of Machine Learning Research* 22 (2021) 3327–3373.
- J. Feng, J. Steinhardt, How do language models bind entities in context?, *arXiv preprint arXiv:2310.17191* (2023).
- K. Meng, D. Bau, A. Andonian, Y. Belinkov, Locating and editing factual associations in gpt, *Advances in Neural Information Processing Systems* 35 (2022) 17359–17372.
- B. Wang, X. Yue, Y. Su, H. Sun, Grokked transformers are implicit reasoners: A mechanistic journey to the edge of generalization, *arXiv preprint arXiv:2405.15071* (2024).

NeurIPS Paper Checklist

1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: Yes, our claims clearly reflect the contributions of this work.

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: We summarize the limitations of this work at the end of the paper.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [NA]

Justification: This work does not present any theoretical results.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [\[Yes\]](#)

Justification: The dataset design, model architecture, and training configurations are described in detail, ensuring that the results can be reliably reproduced.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: We provide the code necessary to reproduce the experiments in the supplementary material.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so No is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: We provide a complete description of the training and testing details in the supplementary material.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [NA]

Justification: Although error bars are not shown, we conducted multiple runs and report averaged results; see Figures 4 and 8.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.

- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: The specific computational resources used in our experiments are described in detail in the supplementary material.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

Answer: [Yes]

Justification: This work meets the requirements of NeurIPS Code of Ethics.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [NA]

Justification: This paper is an experimental study based on observed phenomena and does not have significant societal impact.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.

- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: This work is an experimental study based on observed phenomena, and therefore does not require any safeguards.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: All referenced work has been clearly and appropriately cited throughout the paper.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.

- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. **New assets**

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [NA]

Justification: This work is an experimental and analysis-based study, and therefore does not produce any assets.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. **Crowdsourcing and research with human subjects**

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: This work is an experimental study based on observed phenomena, and therefore does not involve crowdsourcing or human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. **Institutional review board (IRB) approvals or equivalent for research with human subjects**

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: This work is an experimental study based on observed phenomena, and therefore does not involve crowdsourcing or human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

16. **Declaration of LLM usage**

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.

Answer: [Yes]

Justification: This work does not involve anything beyond what is stated in the declaration.

Guidelines:

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (<https://neurips.cc/Conferences/2025/LLM>) for what should or should not be described.

A Experiments compute resources

All experiments were conducted on a server running Ubuntu 22.04.4 LTS. The system is equipped with an Intel Xeon Gold 6133 processor featuring 80 logical cores at 3.0GHz, and 352 GB of RAM. The machine is configured with 8 NVIDIA GeForce RTX 4080 GPUs with 16GB of video memory each.

B Detail of Mamba

The following provides a detailed explanation of the internal computations within the Mamba module, beginning with the specification of the corresponding dimensional representations.

$$\begin{aligned}
 B &:= \text{Batch size,} \\
 S &:= \text{Sequence length,} \\
 D_m &:= \text{dimension of model,} \\
 D_{in} &:= \text{dimension of inner model} \\
 &= \text{expand} \times D_m, \\
 N &:= \text{dimension of hidden state,} \\
 Num_h &:= \text{number of heads,} \\
 H_d &:= \text{dimension of heads,}
 \end{aligned}$$

The computations within the Mamba module can be divided into three components: pre-SSM, SSM, and post-SSM. The following presents a detailed derivation for each submodule.

$$\text{Input: } u \in \mathbb{R}^{(B, S, D_m)},$$

Pre-SSM

$$\begin{aligned}
 zxBCdt &= \mathbf{Linear}_{\text{proj}}(u) \in \mathbb{R}^{(B, S, 2D_{in}+2N+Num_h)}, \\
 (z, xBC, dt) &= zxBCdt \in \mathbb{R}^{(B, S, (D_{in}, D_{in}+2N, Num_h))}, \\
 x_cB_cC_c &= \mathbf{Conv1d}(xBC) \in \mathbb{R}^{(B, S, D_{in}+2N)}, \\
 x_{cf}B_{cf}C_{cf} &= \mathbf{SiLU}(xBC) \in \mathbb{R}^{(B, S, D_{in}+2N)}, \\
 (x_{cf}, B_{cf}, C_{cf}) &= x_{cf}B_{cf}C_{cf} \in \mathbb{R}^{(B, S, (D_{in}, N, N))}, \\
 A_{init} &\sim \mathbf{Uniform}(A_{min}, A_{max}) \in \mathbb{R}^{Num_h}, \\
 A_{log} &= \log(A_{init}) \in \mathbb{R}^{Num_h}, \\
 A &= -\exp(A_{log}) \in \mathbb{R}^{Num_h}, \\
 dtb_{init} &= \exp(\mathbf{rand}(Num_h) \cdot (\log(dt_{max}) - \log(dt_{min})) + \log(dt_{min})) \in \mathbb{R}^{Num_h}, \\
 dtb_{bias} &= dtb_{init} + \log(-\exp(-dtb_{init}) + 1) \in \mathbb{R}^{Num_h}, \\
 dt_{fb} &= \mathbf{softplus}(dt + dtb_{bias}) \in \mathbb{R}^{(B, S, Num_h)}, \\
 \tilde{A} &= A \circ dt_{fb} \in \mathbb{R}^{(B, S, Num_h)}, \\
 \tilde{x} &= x_{cf}.\mathbf{reshape}(B, S, Num_h, H_d) \circ dt_{fb}.\mathbf{reshape}(B, S, Num_h, 1) \in \mathbb{R}^{(B, S, Num_h, H_d)}, \\
 \tilde{B} &= B_{cf}.\mathbf{reshape}(B, S, 1, N) \in \mathbb{R}^{(B, S, 1, N)}, \\
 \tilde{C} &= C_{cf}.\mathbf{reshape}(B, S, 1, N) \in \mathbb{R}^{(B, S, 1, N)}.
 \end{aligned}$$

Here, $\mathbf{Linear}_{\text{proj}}$ denotes a linear transformation, $\mathbf{Conv1d}$ refers to a one-dimensional convolution, and \mathbf{SiLU} represents the SiLU activation function. $\mathbf{Uniform}$ indicates sampling from a uniform distribution, while $\mathbf{softplus}$ denotes the Softplus activation function. \mathbf{rand} refers to drawing a specified number of random values. The \circ symbol denotes element-wise multiplication, and $\mathbf{reshape}$ indicates a dimensional transformation operation.

SSM

$$\begin{aligned}
\widehat{A} &= \text{Mask}_1 \circ \text{repeat}(\widetilde{A}) \in \mathbb{R}^{(B, \text{Num}_h, S, S)}, \\
A_{cumsum} &= \text{cumsum}(\widehat{A}) \in \mathbb{R}^{(B, \text{Num}_h, S, S)}, \\
L &= \exp(\text{Mask}_2 \circ A_{cumsum}) \in \mathbb{R}^{(B, \text{Num}_h, S, S)}, \\
P &= \text{einsum}("BSHN, BSHN \rightarrow BHSS", \widetilde{C}, \widetilde{B}) \in \mathbb{R}^{(B, 1, S, S)}, \\
M &= \text{einsum}("BHSS, BHSS \rightarrow BHSS", L, P) \in \mathbb{R}^{(B, \text{Num}_h, S, S)}, \\
y &= \text{einsum}("BHSS, BSHP \rightarrow BSHP", M, \widetilde{x}) \in \mathbb{R}^{(B, S, \text{Num}_h, H_d)}.
\end{aligned}$$

The Mask_1 operation is used to zero out the diagonal and upper-triangular elements of a matrix, retaining only the elements below the diagonal. The Mask_2 sets the elements above the diagonal to negative infinity. The cumsum operation performs a cumulative sum along the first of the last two dimensions (i.e., across rows) in a matrix of shape (S, S) , summing from top to bottom. The einsum operation denotes the Einstein summation convention, used for concise and flexible tensor contractions.

Post-SSM

$$\begin{aligned}
D &= \mathbf{1}_{D_{in}} \in \mathbb{R}^{D_{in}}, \\
\dot{y} &= y + x_{cf} \cdot \text{reshape}(B, S, \text{Num}_h, H_d) \times D \cdot \text{reshape}(D_{in}, 1) \in \mathbb{R}^{(B, S, \text{Num}_h, H_d)}, \\
\dot{y} &= \dot{y} \cdot \text{reshape}(B, S, D_{in}) \in \mathbb{R}^{(B, S, D_{in})}, \\
y_z &= \dot{y} \cdot \text{SiLU}(z) \in \mathbb{R}^{(B, S, D_{in})}, \\
y_{norm} &= y_z \cdot \frac{1}{\sqrt{\text{mean}(y_z^2, \text{axis} = -1) + \epsilon}} \cdot \mathbf{w} \in \mathbb{R}^{(B, S, D_{in})}, \\
y_{out} &= \text{Linear}_{\text{proj}}(y_{norm}) \in \mathbb{R}^{(B, S, D_m)}.
\end{aligned}$$

Here, $\mathbf{1}$ denotes a vector in which all elements are equal to 1.

$$\text{Output: } y_{out} \in \mathbb{R}^{(B, S, D_m)}.$$

C Data setup

C.1 Composite function task

Standard The total dataset comprises 300,000 samples, and each sequence has a fixed length of 8. Each anchor pair in the training set accounts for 5.6% of the total data, while each anchor pair in the test set constitutes 0.6%. The mapping between anchors and their associated functions is as follows: anchor 1 corresponds to a shift of +5, anchor 2 to +1, anchor 3 to -2, and anchor 4 to -8. During training, all 15 possible anchor pairs are included except for pair 43. Among these, all pairs except 34 are derived from the composition of their corresponding single-anchor functions. Notably, the function for pair 34 is manually set to -6, deviating from the correct compositional result of -10. The test set contains both symmetric and compositional instances of pair 43. The distinction between training and test data is governed by the position of the key token. For sequences of length 8, each position is associated with a congruence class modulo 8. In the training set, the key token is prohibited from appearing in the position whose modulo-8 class matches its own. For example, key 33 cannot appear in the first position, as it belongs to the congruence class $1 \pmod{8}$, which is aligned with the first index. Conversely, in the test set, each key token is required to occupy the position that corresponds exactly to its modulo-8 congruence class. This design ensures that the model has access to the full semantic range of tokens during training while preventing it from relying solely on positional memorization to generalize to the test set.

Full symmetry The total number of datasets is 300,000, and each sequence has a fixed length of 8. Each anchor pair in the training set accounts for 4.5% of the total, while each anchor in the test set accounts for 0.5% of the total. To ensure that the only possible solutions are symmetric,

we simultaneously utilize two sets of correspondences between anchors and functions. This setup prevents the model from simply solving for individual anchor functions and instead forces it to derive symmetric solutions by understanding the symmetry of anchor pairs. There are a total of five anchors: 0, 1, 2, 3, and 4. The 0 anchor is added to balance the data volume between the two function sets, thus avoiding model bias caused by disparities in data quantity. The first set of correspondences between anchors and functions is as follows: 0 corresponds to +2, 1 corresponds to +5, 2 corresponds to +1, 3 corresponds to -2, and 4 corresponds to -8. The second set of correspondences is: 0 corresponds to -9, 1 corresponds to +6, 2 corresponds to -7, and 3 corresponds to +3. All anchor pair functions are composed of either the first or the second set of correspondences. To facilitate the observation of symmetric solutions, anchor pairs 00, 11, 22, 33, and 44 are excluded from this task. Thus, there are a total of 20 anchor pairs in this task. The anchor pairs following the first set of correspondences are 01, 02, 10, 20, 14, 41, 23, 32, 34, 43, a total of 10 pairs. Among them, there are 4 zeros, 4 ones, 4 twos, 4 threes and 4 fours. This design ensures that each anchor has the same amount of data under the two corresponding methods, and the two corresponding methods also have the same amount of data. The anchor pairs following the second set of correspondences are 03, 30, 04, 40, 12, 21, 13, 31, 24, 42. Only 43 is not included in the training set. The purpose of this task is to observe whether the model can discover the symmetry among all anchor pairs and thus output symmetric solutions. The distinction between the training set and the test set in the dataset is based on the position of the key. For sequences of length 8, each position corresponds to a congruence class modulo 8. In the training set, for each key, the key does not fall within the congruence class corresponding to its position. For example, the key 33 cannot appear in the first position of the sequence because it would then be in the congruence class modulo 8 of 1, which corresponds to that position. In contrast, in the test set, each key must exactly fall within the congruence class corresponding to its position. This setup allows the model to learn the meaning of all tokens while preventing it from relying solely on memorization to generalize to the test set.

C.2 Inverse sequence matching task

The total dataset consists of 100,000 samples. The sequence length varies with the number of layers in the Mamba model. Each sequence is generated from a generation set consisting of three distinct elements. It contains five different permutations of the elements in the generation set, where each permutation is followed by a random token that does not belong to the generation set. Each such permutation is referred to as a key sequence.

One of these five sequences is randomly selected and reversed to form the query sequence, which is appended to the end of the original sequence. Between the key sequences and the query sequence, a number of randomly generated tokens equal to the pure convolutional receptive field of Mamba are inserted to prevent the model from retrieving information through convolution. Therefore, the total sequence length can be formally expressed as follows:

$$S = 5 \times (3 + 1) + 3 \times \text{Num}_{\text{layer}} + 3.$$

Here, $\text{Num}_{\text{layer}}$ denotes the number of layers in the Mamba model. For example, in the case of a two-layer Mamba model, the sequence length is 29. To ensure a fair comparison, Mamba and Transformer models with the same number of layers are trained on identical sequence configurations.

The training set constitutes 80% of the total dataset, while the test set and the out-of-distribution (OOD) set each account for 10%. In the training set, the three elements in the generating set do not all belong to the same congruence class modulo 3. In contrast, in the test set, all three elements in the generating set belong to the same congruence class modulo 3. For both the training and test sets, all numerical values in the sequences are drawn from the range 20 to 100. In the OOD set, however, all sequence elements are drawn from the range 101 to 200, and the generation set is not subject to any congruence constraints.

D Training setup

Unless otherwise specified, all tasks in this work adopt the following training parameter settings. The learning rate is initially set to $1e-5$ at the start of training, warmed up to 25 times its initial value within 10 epochs, and then decreases to $1e-5$ via cosine decay at 200 epochs. The training optimizer is AdamW with parameters set as $\beta_1 = 0.9$, $\beta_2 = 0.999$, $\text{eps} = 1e - 8$, and weight decay= $1e-2$.

Meanwhile, gradient clipping is applied with a maximum norm of 1. For composite function tasks, the batch size is set to 2048, while for inverse sequence matching tasks, the batch size is 1024. The loss function used for training is the cross - entropy loss function, which is calculated only for the last token of the model output sequence.

E Experiment detail and result

E.1 Composite function task

E.1.1 Phase diagram of Mamba on the composite function task

Mamba configurations The dimension of model is set to 32, the dimension of hidden state is set to the default value of 128, and the expansion factor is set to the default value of 2, and the activation function employed throughout the model is the Sigmoid Linear Unit (SiLU). To facilitate clear observation, all experiments are conducted using a single head. The convolution kernel length is set to its default value of 4. The values reported at different positions represent the mean results obtained using three fixed random seeds. Across these positions, the only variations are in the model’s depth or initialization.

E.1.2 Experiments with all-one convolution

Mamba configurations To investigate whether the asymmetry of convolution is responsible for Mamba’s difficulty in learning the symmetric solution, we select a model configuration under which Mamba successfully learns the compositional solution but fails to learn the symmetric one. Specifically, we use a five-layer Mamba model with small initialization($\gamma = 1$), while all other settings are consistent with those used in the phase diagram experiments.

E.1.3 Experiments with positional encoding

Mamba configurations To examine whether Mamba’s bias for asymmetry in the composite function task arises from its lack of explicit positional encoding, thus encouraging reliance on convolution. We select a configuration under which Mamba struggles to learn both the composite and symmetric solutions. Specifically, we use a two-layer Mamba model with standard initialization($\gamma = 0.5$). All other settings are kept consistent with those used in the phase diagram experiments.

E.1.4 Experiment under full symmetry

Mamba configurations To validate Mamba’s difficulty in solving the composite function task under fully symmetric settings, we set dimension of model to 128, while keeping all other settings consistent with those used in the phase diagram experiments. The number of layers is varied across 2, 3, 4, 5, and 6, and standard initialization($\gamma = 0.5$) is applied.

Extended results For each configuration, experiments are conducted using three fixed random seeds. The results are shown in the Fig. 10 and Fig. 11. As can be observed, under all configurations, Mamba consistently struggles to solve the composite function task in the fully symmetric setting.

E.1.5 Cosine similarity among convolution weights

Mamba’s asymmetry bias originates from its nonlinear convolution, specifically from the asymmetry of its convolutional kernel parameters. To investigate this asymmetry, we examined the cosine similarity between the convolutional kernels at the beginning and the end of training, as shown in the Fig. 12, and found that they are nearly orthogonal. This indicates that the convolutional parameters exhibit strong asymmetry throughout the training process.

E.2 Inverse sequence matching task

Mamba configurations The dimension of model is set to 128, the dimension of hidden state is set to the default value of 128, and the expansion factor is set to the default value of 2. The activation

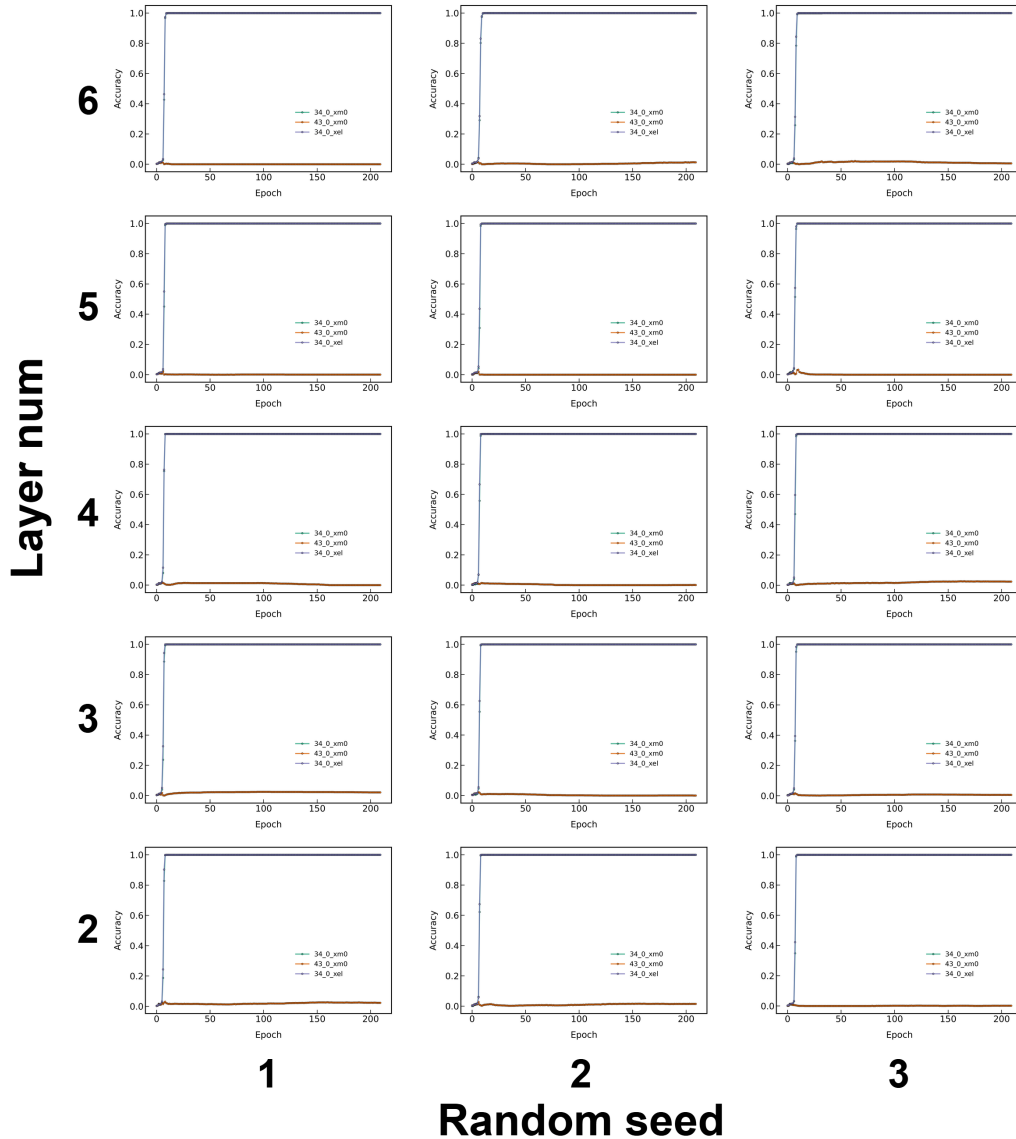


Figure 10: Accuracy of Mamba under fully symmetric setting. The horizontal axis represents the random seed, while the vertical axis corresponds to the number of layers in the Mamba model.

function used is SiLU. To clearly isolate the structural differences between Transformer and Mamba and ensure a fair comparison, the Transformer model also adopts SiLU as its activation function. The convolution kernel length is set to its default value of 4, and the number of heads is fixed at 1.

Transformer configurations The dimension of model is set to 128, consistent with the Mamba configuration. To ensure a fair comparison, the Transformer also adopts a $2\times$ dimensional expansion when generating value vectors, mirroring Mambas setup. Specifically, the dimensions of the query and key vectors are set to 128, while the value vectors have a dimension of 256. Additionally, to ensure a fair comparison and maintain a comparable number of parameters between the two models, the Transformers feedforward network (FNN) uses a hidden dimension of 128, and, like Mamba, it is configured with only a single attention head throughout.

Extended results All experiments are conducted using the same set of random seeds to ensure fairness in comparison.

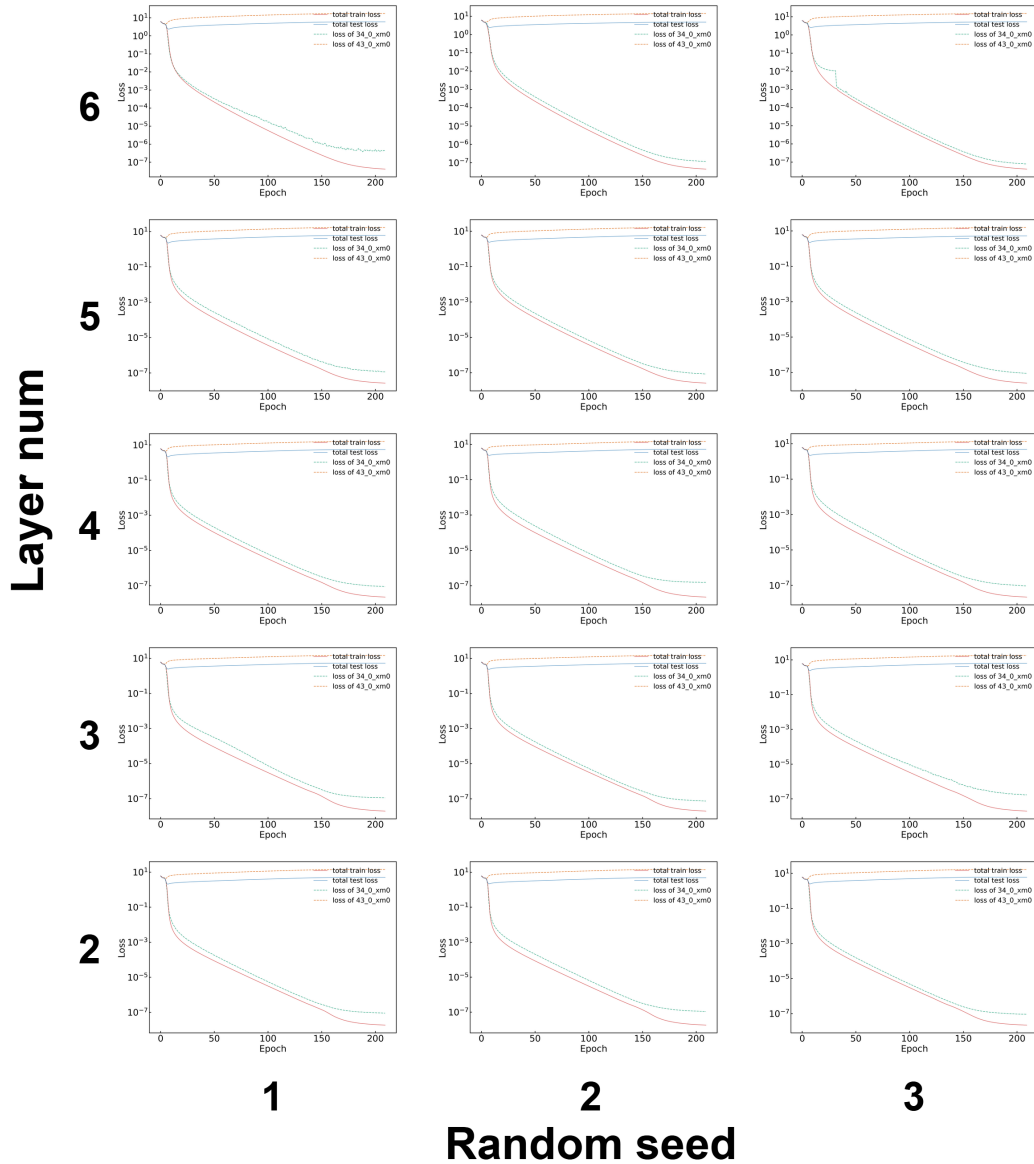


Figure 11: Loss of Mamba under fully symmetric setting. The horizontal axis represents the random seed, while the vertical axis corresponds to the number of layers in the Mamba model.

To evaluate Mamba's difficulty in solving the inverse sequence matching task, we vary the model depth across 2, 3, 4, and 5 layers, and consider different initialization strategies ranging from standard to small initialization ($\gamma = 0.5, 0.6, 0.7, 0.8, 0.9, 1.0$). The detailed results are shown in the Fig. 13 and Fig. 14. It can be observed that Mamba fails to solve the inverse sequence matching task under nearly all configurations.

For the two-layer Transformer, results under different initialization schemes ($\gamma = 0.5, 0.6, 0.7, 0.8, 0.9, 1.0$) are shown in the Fig. 15 and Fig. 16. It can be seen that even with a small number of layers, the Transformer outperforms Mamba.

The results of the modified two-layer Mamba under different initialization schemes ($\gamma = 0.5, 0.6, 0.7, 0.8, 0.9, 1.0$) are shown in the Fig. 15 and Fig. 16. It can be observed that the modified Mamba not only significantly outperforms standard Mamba, but also substantially surpasses the performance of the Transformer.

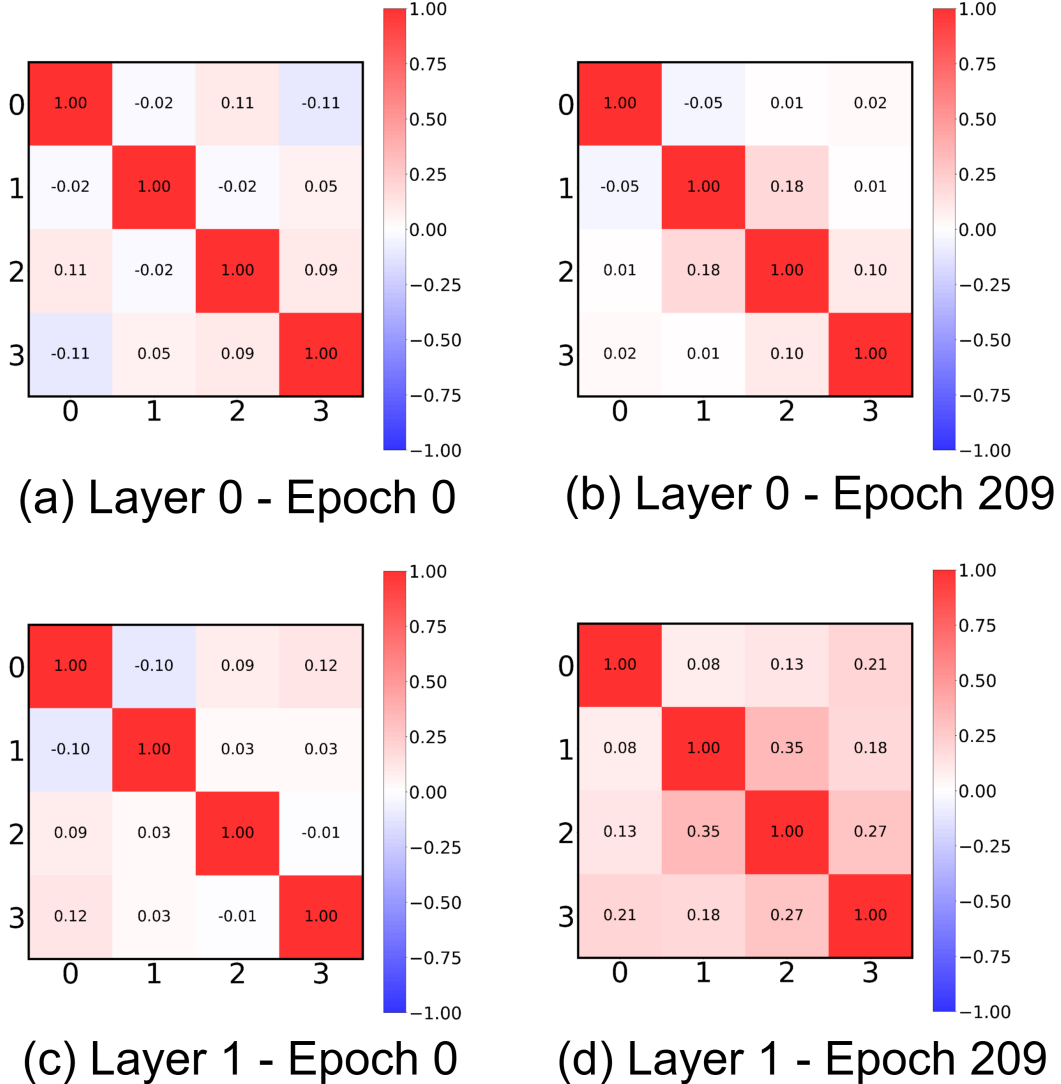


Figure 12: Cosine similarity of convolutional kernel parameters between epoch 0 and epoch 209 for the two-layer Mamba model. The model configuration is consistent with that used in the phase diagram experiments. In this figure, Mamba is configured with two layers and initialized with $\gamma = 1$.

Experiments with alternative architectural modifications From our experiments, we observe that Mamba’s bias toward asymmetry originates from the inherent asymmetry introduced by its nonlinear convolution. However, by introducing a residual connection that bypasses this nonlinear convolution, the impact of such asymmetry can be effectively mitigated. This allows Mamba to leverage the SSM module for information extraction and even achieve performance comparable to or surpassing that of Transformers on the inverse sequence matching task.

It is important to emphasize that, in order for the SSM to extract information in a manner similar to Attention, positional awareness of tokens is essential. Therefore, positional embedding must be incorporated. In fact, the addition of positional embedding plays a significant role in the inverse sequence matching task as well.

In addition to introducing the residual connection that bypasses the nonlinear convolution, we also conducted another set of architectural experiments. Inspired by the original Mamba architecture, another possible way to inject the raw token information is through a "gating mechanism". For example, one might use a pointwise product between the original token and the fused token (after nonlinear convolution) as the input to the SSM. We additionally conducted multiple experiments on

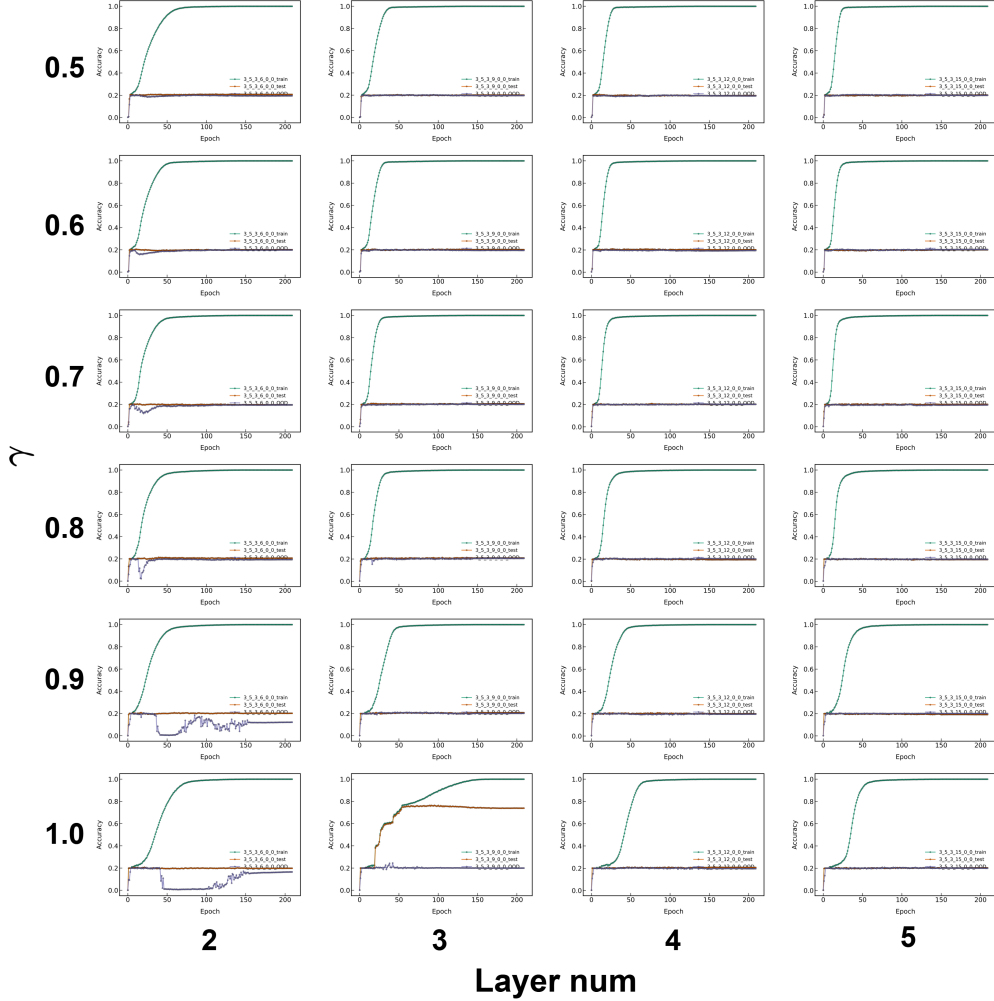


Figure 13: Accuracy of Mamba under different configurations on the inverse sequence matching task. The horizontal axis represents the number of layers in the Mamba model, while the vertical axis corresponds to the initialization scheme.

this approach, and the results are shown in Fig. 17. The following presents the computation process of the gating mechanism. Given the input U to the Mamba block, we have:

$$\begin{aligned}
 (\tilde{U}, Z, dt) &= \text{Linear}(U), \quad U \in R^{(s,d)}, \tilde{U} \in R^{(s,2d+2h)}, Z \in R^{(s,2d)}, dt \in R^{(s,N_h)}, \\
 (B, C, X) &= \sigma(\text{Conv1d}(\tilde{U})), \quad B \in R^{(s,h)}, C \in R^{(s,h)}, X \in R^{(s,2d)}, \\
 (\tilde{B}, \tilde{C}, \tilde{X}) &= (B, C, X) \circ \tilde{U}, \quad \tilde{B} \in R^{(s,h)}, \tilde{C} \in R^{(s,h)}, \tilde{X} \in R^{(s,2d)}.
 \end{aligned} \tag{10}$$

Finally, \tilde{B} , \tilde{C} , and \tilde{X} serve as the input to the SSM.

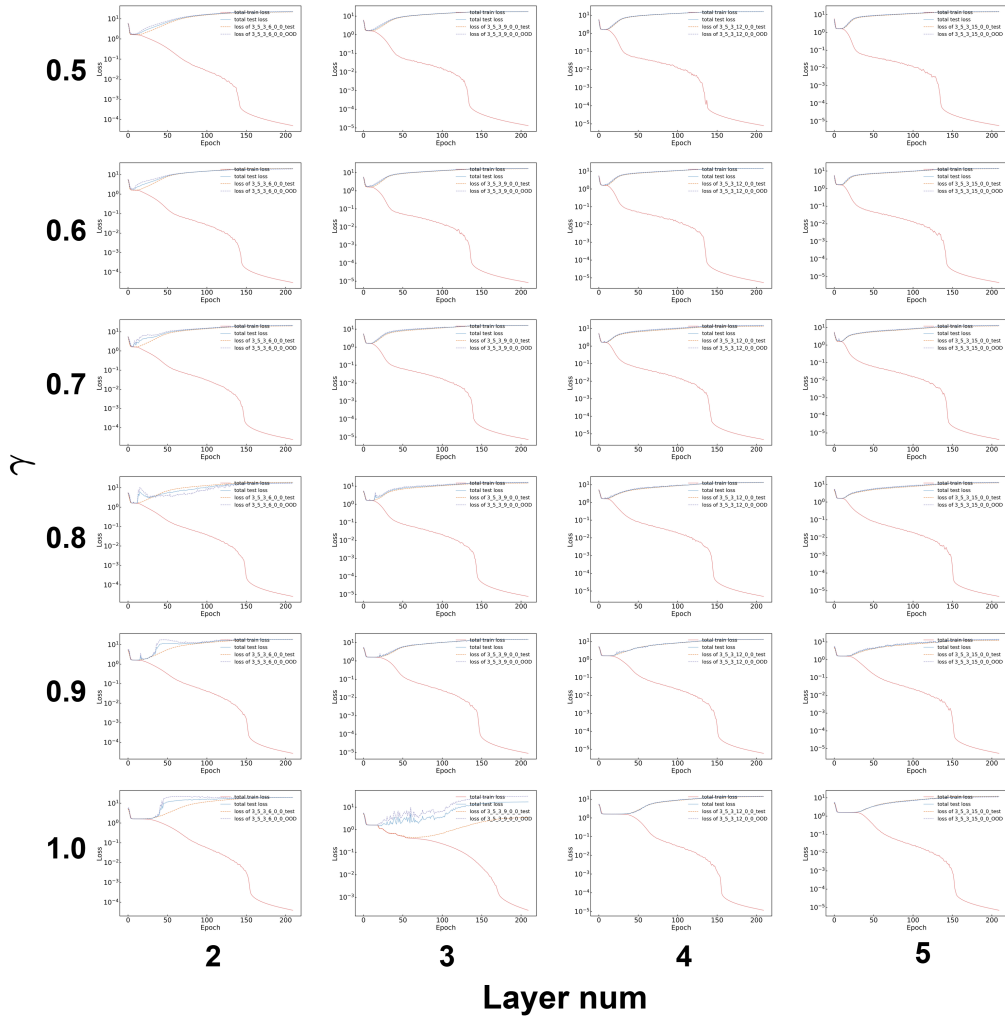


Figure 14: Loss of Mamba under different configurations on the inverse sequence matching task. The horizontal axis represents the number of layers in the Mamba model, while the vertical axis corresponds to the initialization scheme.

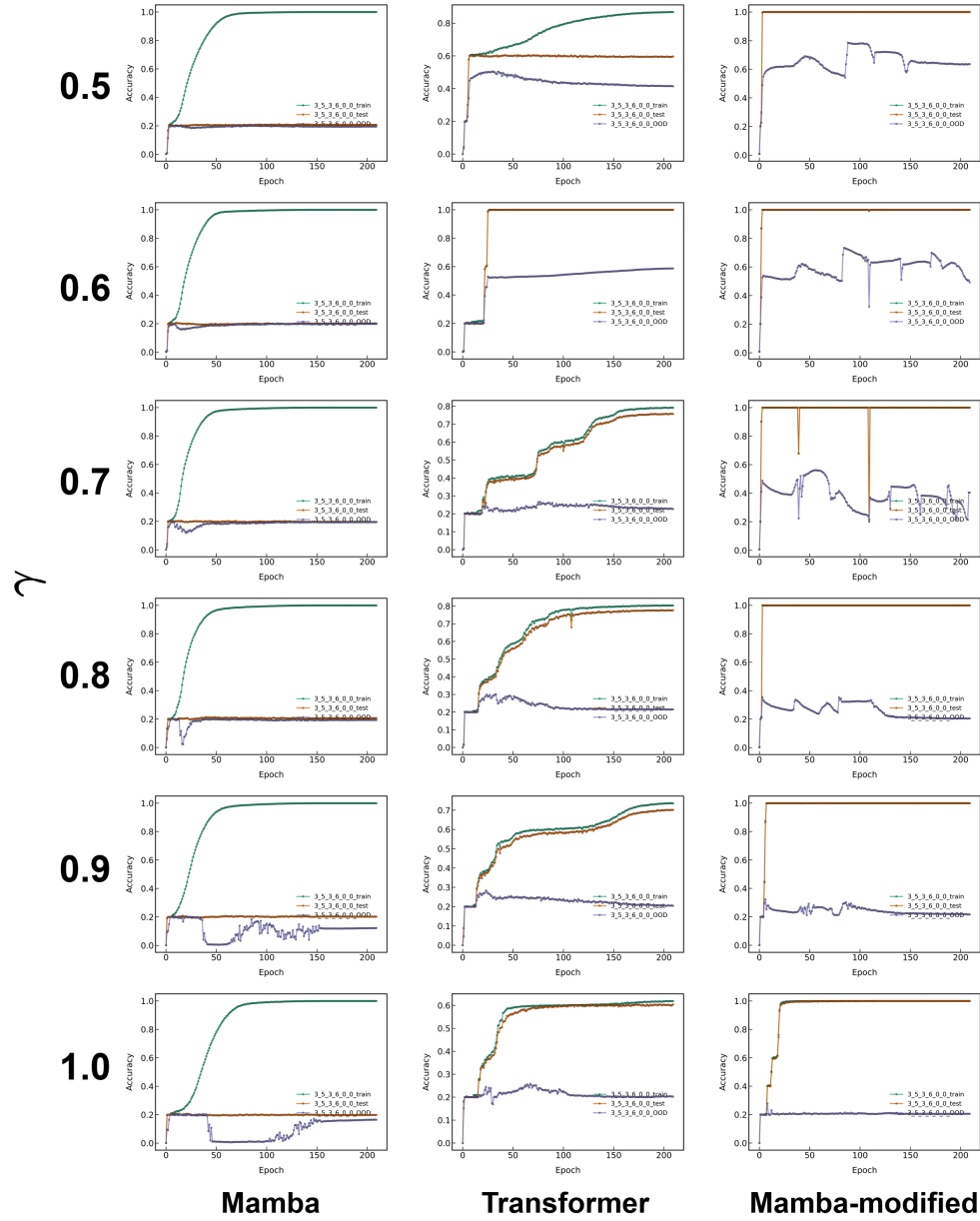


Figure 15: Accuracy of Mamba, Transformer, and modified Mamba on the inverse sequence matching task. Left: Mamba (2-layer), Center: Transformer (2-layer), Right: Modified Mamba (2-layer); all under varying initialization schemes.

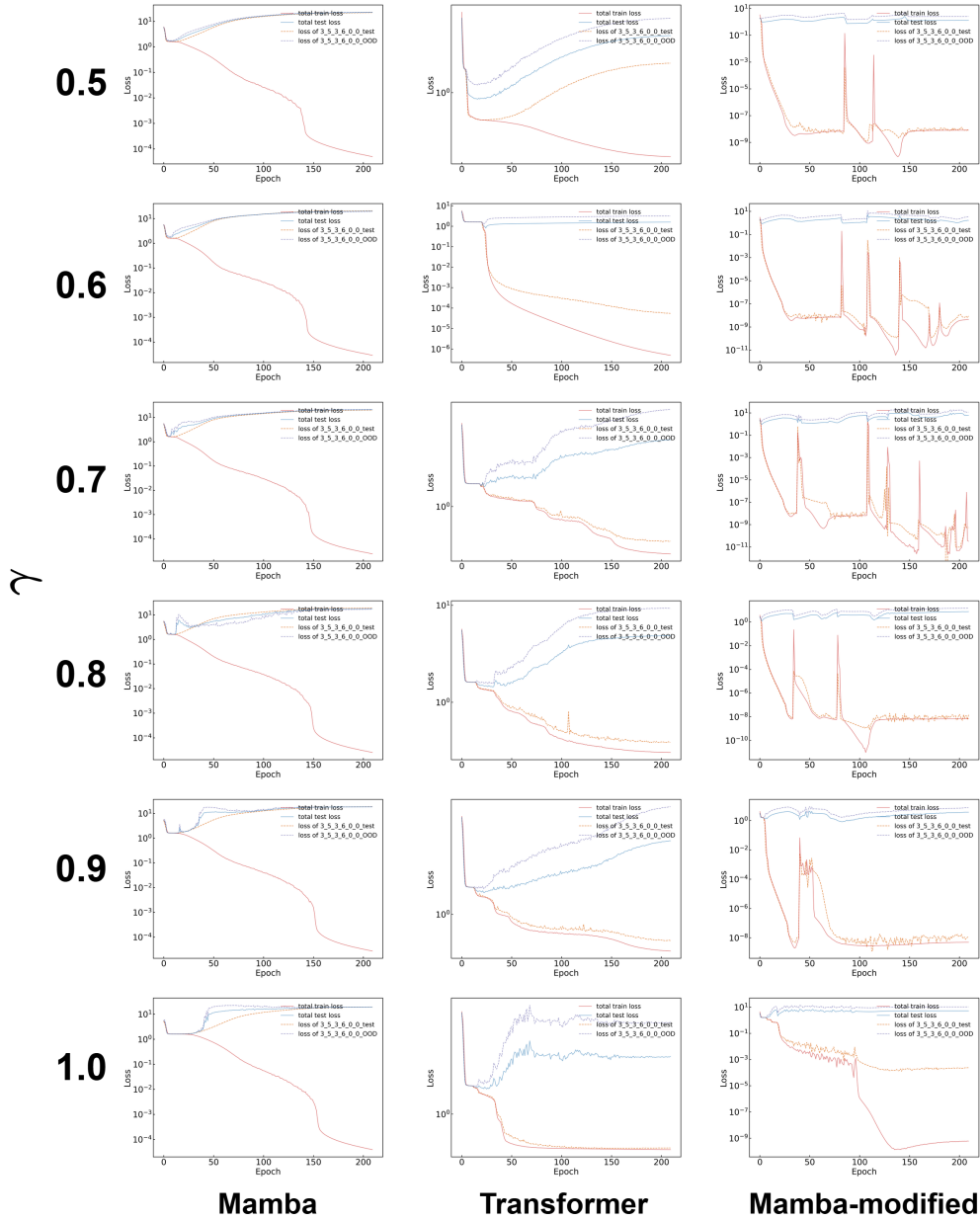


Figure 16: Loss of Mamba, Transformer, and modified Mamba on the inverse sequence matching task. Left: Mamba (2-layer), Center: Transformer (2-layer), Right: Modified Mamba (2-layer); all under varying initialization schemes.

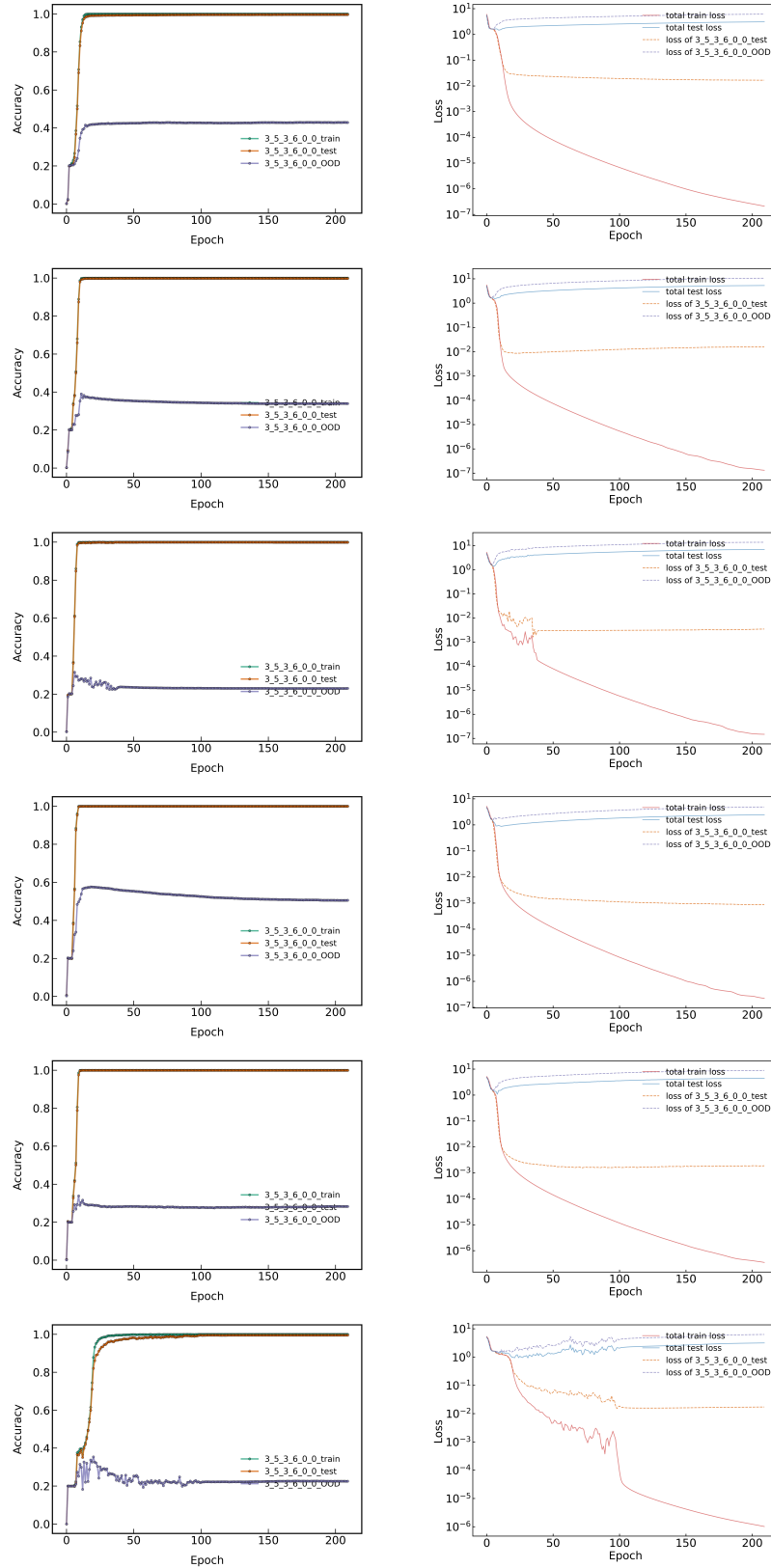


Figure 17: Accuracy (left) and loss (right) curves of Mamba with gate residual connections on the inverse sequence matching task. The architecture of gate-residual Mamba is consistent with that of residual Mamba. From top to bottom, the initialization is set to $\gamma = 0.5$, $\gamma = 0.6$, $\gamma = 0.7$, $\gamma = 0.8$, $\gamma = 0.9$, and $\gamma = 1.0$.