## CAN MAMBA ALWAYS ENJOY THE "FREE LUNCH"?

Anonymous authors

Paper under double-blind review

### ABSTRACT

Transformers have been the cornerstone of current Large Language Models (LLMs); however, its linear growth in overhead during inference with respect to sequence length poses challenges for modeling long sequences. In this context, Mamba has gradually attracted attention due to its constant-level size during inference and existing empirical results have shown that it can perform comparably to Transformers in sequence modeling while offering significant savings. However, one may ask that, can Mamba always enjoy the "free lunch"? In this paper, we focus on analyzing the expressive ability of Mamba from a theoretical standpoint. First, inspired by the connection between Mamba and linear attention, we investigate potential shortcomings of the Mamba when performing the COPY operation. Our results indicate that Mamba with constant size may encounter bottlenecks when handling COPY, while it can achieve perfect performance when the size scales linearly with sequence length. Based on this observation, we analyze Mamba's ability to tackle DP problems when equipped with Chain of Thought (CoT). Our findings suggest that to solve arbitrary DP problems, the total cost of Mamba is comparable to standard and efficient Transformers. However, similar to efficient Transformers, when facing DP problems with favorable properties such as locality, Mamba can provide savings in overhead. Our results contribute to a deeper understanding of Mamba.

026 027 028

029

024

025

000

001 002 003

004

005 006 007

008 009

010

011

012

013

014

015

016

017

018

019

021

### 1 INTRODUCTION

Reccently, Transformer-based large language models (LLMs) have become the mainstream of mod-031 ern neural network architectures due to their outstanding performance across a wide range of tasks (Vaswani et al., 2017; Kenton & Toutanova, 2019; Brown et al., 2020; Dosovitskiy et al., 2020; Min 033 et al., 2022). However, the core component of Transformers—the attention layer—while providing 034 excellent performance, also leads to emerging drawbacks: during training, the computational cost scales quadratically with sequence length, and during inference, the cost scales linearly with sequence length. This limitation becomes increasingly unacceptable when dealing with long sequence 037 tasks. To address this issue, many works have attempted to improve the attention mechanism to 038 reduce its time and storage costs (Tay et al., 2023; Choromanski et al., 2020; Katharopoulos et al., 2020; Beltagy et al., 2020; Child et al., 2019). However, these improved structures often achieve efficiency in the attention layer at the expense of some performance. 040

041 Faced with the scaling challenges of Transformers, the exploration of new model architectures to 042 replace Transformers has gradually come into focus, leading to the development of modern RNN 043 architectures, including RWKV (Peng et al., 2023), RetNet (Sun et al., 2023), and Mamba (Gu & 044 Dao, 2023). Among them, the Mamba architecture (Gu et al., 2021; Gu & Dao, 2023), based on the state space model (SSM), has garnered attention for its performance comparable to Transformers in many sequence modeling tasks Dao & Gu (2024) and vision tasks (Zhu et al., 2024; Xu et al., 2024). 046 These models utilize hardware-aware algorithms during training, resulting in computational costs 047 that scale linearly with sequence length, and require constant-level computation and storage during 048 inference. Mamba's strong performance and computational efficiency make it a strong competitor to Transformers. 050

Despite Mamba demonstrating excellent performance, one can not help but ask, *can Mamba always enjoy such "free lunch", that is, can Mamba always save considerable overhead while maintaining performance of Transformers?* More recent results have revealed Mamba's shortcomings in certain
 tasks, especially those involving retrieval (Arora et al., 2023; Hendrycks et al., 2020; Jelassi et al.,

054 2024). Specifically, Akyürek et al. (2024) study the in-context language learning capabilities of 055 different models and find that Transformers outperformed other models, including Mamba, due to 056 the specialized attention heads. Jelassi et al. (2024) also discover that Transformers are superior to 057 Mamba on tasks that require copying from the input context. Park et al. (2024) point out that Mamba 058 struggles to retrieve vectors from the context of multi-query associative recall (MQAR) (Arora et al., 2023), while Transformers can easily handle it well. Furthermore, Waleffe et al. (2024) conduct experiments on larger models (up to 8B parameters) with a broader range of tasks, discovering that 060 when it comes to in-context learning and recalling information from text, although Mambas can 061 contain the same knowledge as Transformers, it will be more difficult for them to directly copy 062 useful information from history. 063

Although there has been some empirical exploration, the theoretical investigation concerning the 064 above "free lunch" question still remains open to explore. In this paper, we attempt to explore 065 Mamba's expressive ability from a theoretical standpoint. Specifically, inspired by the comparison 066 between Mamba and linear attention mechanism, we first focus on Mamba's ability to perform the 067 COPY operation, which is closely related to the ability to retrieve information from context. Our 068 theoretical results show that Mamba's performance in executing COPY operations is closely related 069 to the size of its model and experiments also confirm the limitations of Mamba in executing copy tasks and retrieval from long sequences. Further, following the setting of Feng et al. (2024); Yang 071 et al. (2024), we explore Mamba's capability to solve DP tasks equipped with Chain of Thought 072 (CoT). We find that to solve arbitrary DP problems, Mamba, standard Transformers and efficient 073 Transformers including Linear and Sparse Transformers (Katharopoulos et al., 2020; Child et al., 074 2019) seem to be on equal footing in terms of inference cost; however, when dealing with *m*-locality 075 DP problems, Mamba may offer significant savings like efficient Transformers. Our results can be concluded as follows: 076

- A constant-sized Mamba may encounter bottlenecks when executing COPY operations (Theorem 1 in Section 4);
- When the size of Mamba scales linearly with the sequence length, it can perform COPY operation accurately (Theorem 2 in Section 4);
- To solve arbitrary DP problems, the total cost required by Mamba is comparable to that of standard and efficient Transformers (Theorem 3 in Section 5);
- When dealing with DP problems that have favorable locality properties, Mamba can bring savings in overhead compared to standard Transformers (Theorem 4 in Section 5).

## 2 RELATED WORK

077

079

081 082

084

085

087

088

SSMs and Attention Mechanism: The attention mechanism is a core component of LLMs (Brown 090 et al., 2020; Touvron et al., 2023). Drawing connections between SSMs and attention is a fascinat-091 ing direction as it not only aids in our understanding of the Mamba structure but also facilitates 092 the transfer of well-established acceleration techniques from attention mechanisms to Mamba (Dao, 2023; Katharopoulos et al., 2020). Based on observations of the similarities between them, Dao 094 & Gu (2024) proposed the state space dual (SSD) layer based on SSMs to achieve significant im-095 provements in training efficiency. Sieber et al. (2024) introduce the Dynamical Systems Framework 096 (DSF), under which attention and SSMs can be directly compared. Additionally, Han et al. (2024) 097 reformulate the structure of SSMs to establish links with linear attention, aiming to investigate the 098 key factors behind success in vision tasks. We follow this convenient reformulation and comparison, based on which we furthermore explore Mamba's ability to perform the COPY operation. 099

100 Comparisons between Transformers and Mamba: More recent works compare the performance 101 of Mamba and Transformers across various tasks from different perspectives. Merrill et al. (2024) 102 theoretically demonstrate that, similar to Transformers, Mamba is also unable to solve state tracking 103 problems such as permutation composition. Jelassi et al. (2024) find that Transformers significantly 104 surpass SSMs when facing tasks related to copying and retrieving information from context. Park 105 et al. (2024) investigate Mamba's capability for in-context learning and demonstrate that Mamba outperforms Transformers in sparse parity learning while it is weaker in tasks involving non-standard 106 retrieval functionality. Similarly, Waleffe et al. (2024) conduct experiments on a larger scale and 107 find that Mamba lag behind Transformers in tasks that require strong copying and long-context

108 reasoning. Our experiments reference the setups of these works and conduct similar investigations. 109 We also note that Jelassi et al. (2024) find theoretical conclusions that appear to be similar to ours, 110 namely that generalized SSMs (GSSMs) including Mamba cannot copy uniformly input sequences 111 unless the size of the state space grows linearly with the sequence length. Although our work 112 also focuses on "copy", Jelassi et al. (2024)'s definition of the copy task is at the sequence level, investigating the ability of GSSMs to replicate entire sequences that satisfy some distribution and 113 providing a lower bound for their state space memory. In contrast, our COPY operation is defined 114 at the token level, exploring the conditions under which Mamba can copy some specific historical 115 token and providing constructions for Mamba with linear-scaling size to achieving this operation. 116

117 Transformers and modern RNNs with CoT: Chain-of-Thought (CoT) (Wei et al., 2022) is em-118 ployed to enhance the performance of LLMs by enabling them to provide step-by-step reasoning before arriving at a final answer. It has been shown theoretically that Transformers with CoT exhibit 119 significantly improved expressive power, allowing them to solve more complex problems compared 120 to Transformers without CoT (Merrill & Sabharwal, 2023b; Feng et al., 2024; Merrill & Sabharwal, 121 2023a; Li et al., 2024; Yang et al., 2024). Our analysis of Mamba equipped with CoT follows the 122 framework set by Feng et al. (2024); Yang et al. (2024) in their analysis of dynamic programming 123 (DP) problems. Furthermore, Wen et al. (2024) examine the effect of CoT on enhancing the expres-124 sive power of modern RNNs including SSMs by drawing connections with a Turing machine with 125 O(log(n)) space. Their findings indicate that when using CoT, RNNs with log(n) bit memory will 126 have strictly stronger representation power than those without CoT. Different from this, our work 127 explores the ability of Mamba equipped with CoT from the perspective of solving DP problems fol-128 lowing the setting of Feng et al. (2024) and show constructions for Mamba layers with linear-scaling 129 size relative to the sequence length to solve DP problems.

130 131

132

### 3 PRELIMINARIES

In this section, we introduce the Mamba structure that we focus on and its reformulated form firstly introduced by Han et al. (2024), which facilitates a better understanding of the connection between Mamba and linear attention as illustrated in Section 4.1. It should be noted that to better distinguish different types of variables, in this paper, we use bold uppercase letters to represent matrices such as *A*, bold lowercase letters to represent vectors such as *a*, and all non-bold letters to represent scalars such as *a* and  $\Delta$ . This may differ slightly from the notations used in some Mamba-related literatures (Dao & Gu, 2024; Zhu et al., 2024), where uppercase letters are used to describe *A*, *B*, *C* in SSMs.

140 State Space Model: The state space model (SSM) is inspired by the continuous system that maps a 141 scalar input  $x(t) \in \mathbb{R}$  to its output  $y(t) \in \mathbb{R}$  through a high-dimensional hidden state  $h \in \mathbb{R}^{d_h}$  (Gu 142 & Dao, 2023; Dao & Gu, 2024; Han et al., 2024; Zhu et al., 2024; Han et al., 2024). Specifically, 143 this system can be written as:

$$\boldsymbol{h}'(t) = \boldsymbol{A}\boldsymbol{h}(t) + \boldsymbol{b}\boldsymbol{x}(t), \quad \boldsymbol{y}(t) = \boldsymbol{c}^T\boldsymbol{h}(t) + d\boldsymbol{x}(t),$$

where  $A \in \mathbb{R}^{d_h \times d_h}$  denotes the evolution parameters,  $b, c \in \mathbb{R}^{d_h}$  are projection parameters and dis a scalar parameter. The above continuous system can be discretized using transformation called zero-order hold (ZOH), resulting in a discrete version that can be used for neural networks. In this process, A, b will be transformed as  $\overline{A}, \overline{b}$ . The discrete version for SSM can be written as:

149 150

144

$$\boldsymbol{h}_i = \overline{\boldsymbol{A}} \boldsymbol{h}_{i-1} + \overline{\boldsymbol{b}} x_i, \quad y_i = \boldsymbol{c}^T \boldsymbol{h}_i + dx_i,$$

where  $\overline{A} = \exp(\Delta A)$ ,  $\overline{b} = (\Delta A)^{-1}(\exp(\Delta A) - I) \cdot \Delta b \approx \Delta b$  and  $\Delta \in \mathbb{R}$  is a timescale parameter. The matrix A is typically assumed to have certain structures such as being diagonal, leading to the structured SSMs (Gu et al., 2022; Gupta et al., 2022).

154 Selective State Space Module: To enhance the SSM, Mamba makes the parameters  $b_i, c_i, \Delta_i$ 155 dependent on different inputs  $x_i$ . More specifically, A is set to be diagonal resulting in that 156  $\overline{A}_i h_{i-1} = \tilde{a}_i \odot h_{i-1}$  where  $\tilde{a}_i = \exp(\Delta_i a)$ ,  $a = \operatorname{diag}(A)$  and  $\odot$  denotes the element-wise 157 product. In addition,  $\overline{b}_i x_i = \Delta_i b_i x_i = b_i (\Delta_i \odot x_i)$ . Thus, this transformation ultimately results in: 158  $b_i = \tilde{a}_i \odot b_i + b_i (\Delta_i \odot x_i) = a_i^T b_i + d \odot x_i$ 

$$\boldsymbol{h}_i = \tilde{\boldsymbol{a}}_i \odot \boldsymbol{h}_{i-1} + \boldsymbol{b}_i (\Delta_i \odot x_i), \quad y_i = \boldsymbol{c}_i^T \boldsymbol{h}_i + d \odot x_i.$$

Furthermore, to extend the case of processing scalar inputs  $x_i$  to vectors  $x_i \in \mathbb{R}^d$ , Mamba performs the above operations on each dimension independently, which can be formalized as:

$$\boldsymbol{H}_{i} = \boldsymbol{A}_{i} \odot \boldsymbol{H}_{i-1} + \boldsymbol{b}_{i} (\boldsymbol{\Delta}_{i} \odot \boldsymbol{x}_{i})^{T}, \quad \boldsymbol{y}_{i} = \boldsymbol{H}_{i}^{T} \boldsymbol{c}_{i} + \boldsymbol{d} \odot \boldsymbol{x}_{i}, \tag{1}$$

where  $\widetilde{A}_i = [\widetilde{a}_i^{(j)}]_{j=1}^d \in \mathbb{R}^{d_h \times d}$ ,  $b_i = W_b x_i \in \mathbb{R}^{d_h}$ ,  $c_i = W_c x_i \in \mathbb{R}^{d_h}$ ,  $\Delta_i =$ Softplus $(W_{\Delta}^2 W_{\Delta}^1 x_i) \in \mathbb{R}^d$  and  $W_{\Delta}^1, W_{\Delta}^2$  are linear projection parameters. Thus, given the input  $X = [x_i]_{i=1}^N \in \mathbb{R}^{d \times N}$ , we denote the output of the SSM module in Mamba as Y =SSM(X)162 163 164 165 where  $Y = [y_i]_{i=1}^N$  and  $y_i$  follows Eq (1). This formalization was introduced by Han et al. (2024) 166 to build a bridge between Mamba and linear attention and here we follow this form. 167

**Mamba Layer:** Given some input sequence  $X = [x_i]_{i=1}^N \in \mathbb{R}^{d \times N}$ , the input will be processed through stacked Mamba layers where each layer can be viewed as consisting of a residual connection and a Mamba block  $f^{(l)} : \mathbb{R}^d \to \mathbb{R}^d$ . Specifically, this process can be formulated as 168 169 170

$$\boldsymbol{X}_{l} = \boldsymbol{X}_{l-1} + \boldsymbol{f}^{(l)}(\boldsymbol{X}_{l-1}), \quad l = 1, 2, \dots, L$$

$$(2)$$

$$(2)$$

$$(2)$$

$$(3)$$

(2)

 $\boldsymbol{f}^{(l)}(\boldsymbol{X}_{l-1}) = \boldsymbol{W}_{3}^{(l)} \cdot \text{SSM}(\boldsymbol{W}_{1}^{(l)}\boldsymbol{X}_{l-1} + \boldsymbol{b}_{1}^{(l)}) \odot \sigma(\boldsymbol{W}_{2}^{(l)}\boldsymbol{X}_{l-1} + \boldsymbol{b}_{2}^{(l)}), \quad (3)$ where  $\sigma(\cdot)$  denotes SiLU activation function. A Mamba block combines the output of the SSM 173 174 module with a gated MLP (Gu & Dao, 2023; Chowdhery et al., 2023; Shazeer, 2020), which can be 175 viewed as consisting of two branches. Here we call the branch with the SSM module as "the SSM 176 branch" and the other as "the gated branch". It should be noticed that for the sake of simplifying the 177 analysis, we focus on the core components of Mamba including SSM module and the gated MLP 178 while ignoring other structures including the  $\sigma(\cdot)$  before SSM module, Layer Normalization and 1D 179 Convolutions. The simplified Mamba structure that we focus on can be illustrated by Figure 1. In addition, to avoid confusion, we clarify that the SSM module here is specifically the Selective State Space Module used in Mamba and this will be consistent throughout the subsequent sections. 181

182 183

#### CAN MAMBA ALWAYS PERFORM COPY PERFECTLY? 4

185 In this section, we firstly interpret the reformulated SSM module introduced in Section 3 as a spe-186 cial linear attention. Then based on this observation, we explore the capability of Mamba to execute COPY operations during inference, which is highly related to the ability to perform in-context learn-187 ing and retrieve information. 188

189 190

191

196 197

205 206 207

### 4.1 VIEWING MAMBA AS LINEAR ATTENTION

The attention mechanism is the key to the success of the Transformer architecture. Recent works has 192 explored the relationship between Mamba and attention mechanisms, particularly linear attention 193 from different perspectives Han et al. (2024); Dao & Gu (2024); Sieber et al. (2024). The linear 194 causal attention mechanism can be formalized as: 195

$$\boldsymbol{y}_{i} = \sum_{j=1}^{i} \boldsymbol{v}_{j} \boldsymbol{k}_{j}^{T} \boldsymbol{q}_{i} = \sum_{j=1}^{i} (\boldsymbol{q}_{i}^{T} \boldsymbol{k}_{j}) \boldsymbol{v}_{j} = \sum_{j=1}^{i} a_{ij} \boldsymbol{v}_{j}, \qquad (4)$$

where  $q_i$ ,  $k_i$ ,  $v_i$  are usually interpreted as query, key, value respectively and  $a_{ij}$  denotes the attention scores of the *i*-th token to the *j*-th one. In attention mechanisms in Transformers, there exists  $a_{ij} > 0$ 200 for all  $j \leq i$  and  $\sum_{i=1}^{j} a_{ij} = 1$ , which can be implemented by Softmax function, or approximated by kernel methods (Katharopoulos et al., 2020; Choromanski et al., 2020). 201 202

On the other hand, given the input sequence  $[x_i]_{i=1}^N$  and recalling Eq (1), that the output of the SSM 203 module will have the following form when we set  $H_0 = O$  and d = 0: 204

$$\boldsymbol{y}_{i} = (\boldsymbol{\Delta}_{i} \odot \boldsymbol{x}_{i})\boldsymbol{b}_{i}^{T}\boldsymbol{c}_{i} + \sum_{j=1}^{i-1} \left[\boldsymbol{\Pi}_{j} \odot (\boldsymbol{\Delta}_{j} \odot \boldsymbol{x}_{j})\boldsymbol{b}_{j}^{T}\right]\boldsymbol{c}_{i},$$
(5)

where  $\Pi_j = \widetilde{A}_i \odot \widetilde{A}_{i-1} \odot \cdots \odot \widetilde{A}_{j+1}$ . We notice that since in practice all elements of  $\Delta$  are positive and A is set to be negative (Gu & Dao, 2023; Dao & Gu, 2024; Han et al., 2024), so that 208 209 the elements of  $\hat{A}_i$  in Eq (1) belong to the interval [0, 1]. For the sake of simplicity in analysis, we 210 replace the matrix  $\hat{A}_i$  with a constant  $a_i$  (i.e., considering the case where all elements of  $\hat{A}_i$  are the 211 same), where  $a_i \in [0, 1]$  (Dao & Gu, 2024). In fact, the subsequent analysis can be easily extended 212 to the normal case where the elements of matrix  $A_i$  are different. Then, Eq (5) can be rewritten as 213

214  
215 
$$\boldsymbol{y}_i = \sum_{j=1}^i \alpha_j (\boldsymbol{\Delta}_j \odot \boldsymbol{x}_j) \boldsymbol{b}_j^T \boldsymbol{c}_i = \sum_{j=1}^i \alpha_j (\boldsymbol{c}_i^T \boldsymbol{b}_j) (\boldsymbol{\Delta}_j \odot \boldsymbol{x}_j), \quad (6)$$



Figure 1: The illustration of the simplified Mamba layer we focus on. Left Part: A Mamba layer can be composed of a Mamba block with the residual connection; The Mamba block uses a gated MLP to control the output of the SSM module, where we call the branch with the SSM module as "the SSM branch" while the other as "the gated branch"; **Right Part:** The SSM module used in Mamba can be rewritten in a form similar to linear attention, where  $\Delta_i$ ,  $b_i$ , and  $c_i$  in SSM are all derived from the current  $x_i$ , similar to  $v_i$ ,  $k_i$ , and  $q_i$  in linear attention respectively.

where  $\alpha_j = \prod_{k=j+1}^i a_k$  for  $j \le i-1$  and  $\alpha_i = 1$ . In this form, we can observe that it bears similarities to linear attention without normalization in Eq (4), where  $(\Delta_j \odot x_j)$ ,  $b_j$ ,  $c_i$  corresponds to  $v_j$ ,  $k_j$  and  $q_i$  respectively and  $c_i^T b_j$  acts like attention scores  $a_{ij}$ . Considering  $\alpha_{j-1} \ge \alpha_j$  and  $\alpha_j \in [0, 1]$  for all  $j \le i$ , the main difference is that each term in Eq (6) is weighted by a coefficient  $\alpha_j$  to achieve the forgetting of inputs at longer distances while the attention mechanism uses the constraints for attention scores imposed by Softmax function to make sure the scaling of outputs.

### 4.2 THE TRADE-OFF OF MAMBA WHEN PERFORMING COPY

Based on the observation of the connection between the SSM module and attention mechanism, we investigate the capability of Mamba to recover historical inputs, which is foundational for the model to process information based on context. The COPY operation we focus on is defined as follows:

The *L*-local matching set  $S_i$  describes the indices of historical keys  $b_j$  that is highly relevant to the current query  $c_i$  within a local window of length *L*, that is, the "attention scores  $|c_i^T b_j|$ " is lower-bounded by  $\delta$ . Thus, the position of the historical record we most want to replicate should be within this highly relevant set. It should be noted that here we do not specify particular conditions for the positions we want to copy for a given  $i \in [N]$ , which are usually determined by specific tasks. Therefore, our analysis can be applied to general scenarios. We then make the following assumption:

### 260

236

243

244

**Assumption 1.** For the given sequence  $X = [x_i]_{i=1}^N$ , the following conditions holds:

261 262

263

- For  $i \in [N]$ ,  $\boldsymbol{c}_i^T \boldsymbol{b}_{nos(i)} \ge \rho$  where  $\rho \ge \delta$  and  $0 \le |\boldsymbol{c}_i^T \boldsymbol{b}_i| < \delta$  for all  $j \notin S_i$ .
- For any  $i \in [N]$ , there exists  $\sum_{j \in S_i} \alpha_j |\mathbf{c}_i^T \mathbf{b}_j| < 1$ .
- The first condition of Assumption 1 places a restriction on the attention score between the historical keys and the current query: for those positions not within the *L*-local matching set  $S_i$ , the relevance between the two will be strictly constrained within  $\delta$ ; whereas for the position we want to replicate, its relevance is lower-bounded by  $\rho$ , ensuring a certain gap from other historical records. The second condition of Assumption 1 imposes a constraint on the stability of the output caused by the records in  $S_i$ , that is,  $\|\sum_{i \in S_i} \alpha_j c_i^T b_j v_j\| < 1$  when  $\|v_j\| \leq 1$  for  $j \in S_i$ .

270 Now we present the following result:271

**Theorem 1** (Perform COPY operation with constant size). Under Assumption 1, given a input sequence  $x_1, x_2, ..., x_N$  and for any  $\epsilon > 0$ , there exists a Mamba block with constant size that can approximate the COPY operation, that is, for  $i \in [N]$ , we have  $||y_i - o_i|| \le \epsilon$  if the following condition is satisfied:

275 276 277

278

$$\rho \ge \left(1 - \frac{\epsilon}{2M \|\mathbf{\Delta}\|_{\infty}}\right) \frac{1}{\alpha_{pos(i)}} + \frac{\delta}{2} \left(\frac{a_{\max,<}}{1 - a_{\max,<}} + \frac{\left(\frac{1}{a_{\min,>}}\right)^{L-1} - 1}{1 - a_{\min,>}}\right),\tag{7}$$

where  $a_{\max,<} = \max_{1 \le j < pos(i)} a_j$  and  $a_{\min,>} = \min_{pos(i) < j \le i} a_j$ . Moreover, in such cases, there exists  $\delta \le \frac{\epsilon}{(L-1)M \|\Delta\|_{\infty}}$  when  $L \ge 2$ .

The proof of Theorem 1 can be seen in Appendix A.1. Theorem 1 presents the trade-off between the historical coefficient  $a_j$ , the window length L of  $S_i$ , and the lower bound of the attention score  $c_i^T b_j$ when performing the COPY operation with constant size. At first glance, the lower bound  $\rho$  exhibits exponential growth with respect to L, that is,  $\rho = \Omega \left( (1/a_{\min,>})^L \right)$ , which seems unacceptable.

286 In addition, when  $a_{\min,>}$  and L are smaller and  $a_{\max,<}$  is larger, the aforementioned condition (7) 287 will be satisfied more easily. More specifically, as seen from (7), when we want to copy the vector 288  $v_{pos(i)}$  from the hidden state, the terms before pos(i) should be forgotten sufficiently  $(a_{\max,<}$  is 289 smaller) so that the record remembered by the hidden state is as pure as possible. In addition, the 290 forgetting coefficient after pos(i) should not be too small (correspondingly,  $a_{\min,>}$  should be as large as possible) to make sure that  $v_{pos(i)}$  is recorded to some extent. Futhermore, the maximum 291 distance L between the copied position and the current position should not be too far. In these cases, 292 the lower bound of the relevance score  $c_i^T b_{pos(i)}$  will be smaller, which means that it will be easier to 293 achieve the COPY operation. Moreover, it should be noticed that as the possible position of pos(i)move further away from the current position (L grows), the attention scores of irrelevant historical 295 records will decrease, that is,  $|c_i^T b_j|$  will be smaller for  $j \notin S_i$ . 296

One cannot help but think that the sufficient condition for Mamba to achieve the COPY operation illustrated by Theorem 1 is extremely stringent as the lower bound of the attention score for the historical record we want to copy grows exponentially with *L*. This naturally leads to the question: *will Mamba perform the COPY operation more easily when we increase the model size*? We point out that when the size of Mamba increases linearly with the length of the input sequence, Mamba will be capable of accurately restoring the historical records. Below, we present our results:

**Theorem 2** (Perform COPY operation with linear-scaling size). Given sequence  $x_1, x_2, ..., x_N \in [-M, M]^d$ , there exists a Mamba block with size O(N) that can perform the defined COPY operation, that is,  $y_i = o_i$  for any  $i \in [N]$ . Moreover, the  $l_{\infty}$  norm of the Mamba block parameters is upper bounded by O(ploy(M, N)).

307 The proof of Theorem 2 can be found in Appendix A.2. Theorem 2 is based on a simple intuition: 308 when the size grows linearly with the length of input sequence, the model will have enough space to 309 store these historical records and therefore can retrieve them. All parameters being upper bounded by O(ploy(M, N)) means that the problem can also be solved by the same Mamba block with log(N)310 precision, which has been also adopted in previous works (Merrill & Sabharwal, 2023b;a; Feng 311 et al., 2024; Yang et al., 2024; Wen et al., 2024). Furthermore, here we only provide an existence 312 construction, and whether the Mamba layers will actually learn these constructions is beyond the 313 scope of our analysis. It should be noticed that during inference, such a Mamba block will have the 314 same cost as Transformers at each step(growing linearly with length). Based on this observation, we 315 will elaborate in Section 5 that when faced with the DP problems, Mamba will incur the same order 316 of overhead as the Transformer.

- 317
- 318 319

4.3 MAMBA EMPIRICALLY WEAKER IN COPY TASKS THAN TRANSFORMERS

Although theoretical analysis shows that Mamba might face difficulties in performing COPY opera tion, in practice, Mamba may mitigate this issue through multi-layer stacking, and the actual number
 of parameters might be sufficient to handle tasks of a certain scale. Therefore, following previous
 works (Jelassi et al., 2024; Waleffe et al., 2024), we conduct experiments on both synthetic data and
 more realistic task to explore the practical performance of Mamba.



Figure 2: Experiment on the copy task and more realistic phonebook task. Left Part: The training process of Transformers and Mamba models of different sizes on the copy task. Center Part: The performance of Mamba when the length of the string to be copied is changed. **Right Part:** The performance of pre-trained Mamba as the length of the phone book increases. 335

#### **EXPERIMENTS ON SYNTHETIC DATA** 4.3.1

338 First, we conduct experiments on the copy task introduced by Jelassi et al. (2024). In this task, 339 models are given a string of a certain length and then are required to copy this string exactly as out-340 put. During training, we uniformly sample a length from the interval  $[N_{\min}, N_{\max}]$  and then sample 341 characters from the alphabet at each position to obtain a string as a instance. During evaluation, we sample strings of fixed length  $N_{\rm max}$  as inputs, and then we ask models to repeat them while count-342 ing the ratio of correctly copied characters as accuracy. More experimental details can be found in 343 Appendix A.5. 344

345 For the left part of Figure 2, we firstly compare Transformer and Mamba of similar size (126M 346 and 135M, labeled as TF and Mamba) while fixing the string length. The result shows that while 347 both can perform the copy task finally, Transformer learns the task quickly whereas Mamba requires nearly seven times the number of examples. For smaller models, on the one hand, we find that the 348 performance of Transformer is hardly affected when we halve the number of layers (63M, marked 349 as TF-small); on the other hand, we consider two approaches for Mamba: (i) reducing the number 350 layers (67M, labeled as Mamba-small-L), which causes Mamba to require more examples to achieve 351 similar accuracy; (ii) reducing the hidden size (67M, labeled as Mamba-small-D), in such case 352 Mamba learn the copy task even more slowly; (iii) reducing the hidden size while increasing 353 **layers** (69M, labeled as Mamba-small-LD), in which case Mamba can not learn the task with  $5.8 \times$ 354  $10^5$  examples and the training process becomes unstable. Our results indicate that for the copy task, 355 the hidden size of Mamba seems more important than the number of layers at the same model size, 356 while Transformers are consistently better than Mamba.

357 For the center part of Figure 2, we change the maximum length  $N_{\rm max}$  while maintaining the model 358 size. The results showed that when  $N_{max}$  are 10, 20, and 30, the number of training examples needed 359 for Mamba to learn the task (achieving 95% accuracy in 3 consecutive evaluations) is approximately 360  $1.36 \times 10^4$ ,  $3.52 \times 10^4$ , and  $1.15 \times 10^5$  respectively, indicating that the required number grows 361 more than linearly relative to  $N_{\text{max}}$ . Notably, when  $N_{\text{max}} = 40$ , Mamba will be unable to learn the 362 task within  $5.8 \times 10^5$  examples. Additionally, we observe instability in Mamba's training as  $N_{\rm max}$ increases. In contrast, Transformer can still learn quickly and maintain stability even at  $N_{\text{max}} = 40$ , 364 which again indicate that Transformer outperforms Mamba in executing copy operations.

365 366

367

332

333

334

336

337

#### 4.3.2 **EXPERIMENTS ON PHONEBOOK TASK**

368 For more realistic tasks, we consider the "phonebook" task following Waleffe et al. (2024); Jelassi 369 et al. (2024). In this task, models are given a phone book consisting of N names and their corresponding phone numbers, and then models are asked in a few-shot manner to provide the phone 370 number for some given person in the phone book, for example, "Bob: 111111; Alice: 222222; Tom: 371 { } ", which relates to the ability to copy at specified positions. We examine the pre-trained Mamba 372 models of sizes 370M, 1.4B, and 2.8B (Gu & Dao, 2023). 373

374 It can be seen from the right part of Figure 2 that when the length N of the phone book increases, the performance of models declines to a certain extent while larger models are better at retrieving 375 the required information. Additionally, we notice that there will be some improvement in perfor-376 mance when we provide Mamba with task-related information in advance, such as adding the prompt 377 "Please remember the phone number of **Tom**" at the beginning (labeled as Mamba(2.8B)++). This

suggests that informing the model about the task in advance helps it to perform better in selective
 memory during processing subsequent inputs, thereby enhancing the effectiveness of the copy.

381 382

384

385

386

387

388

389

390

391

392

393 394

395

396

397

414

419

420 421

422

423 424

425

426

427

428

## 5 THE EXPRESSIVE POWER OF MAMBA EQUIPPED WITH COT

Although Mamba may face certain bottlenecks when handling copy tasks, another interesting question is: when augmented with other techniques, such as Chain of Thought (CoT), will Mamba see an improvement in its capabilities? Now, we turn our attention to the expressive ability of Mamba equipped with CoT to solve dynamic programming (DP) problems. Specifically, following the setup by Feng et al. (2024); Yang et al. (2024), a DP problem can be described by input sequences  $\{s^{(1)}, s^{(2)}, \ldots, s^{(N)}\}$ , state space  $\mathcal{I}$ , transition function  $f_{\mathcal{T}}$ , and aggregation function  $f_{\mathcal{A}}$ . Each of component can be described as follows:

• Input sequences: We use  $\{s^{(1)}, s^{(2)}, \dots, s^{(N)}\}$  to denote the input of the sequences and the vector  $\boldsymbol{n} = [|s^{(1)}|, |s^{(2)}|, \dots, |s^{(N)}|]^T$  to describe the scale of the problem, where  $|s^{(i)}|$  denotes the length of the *i*-th sequence.

• State Space: For a given DP problem, the state space  $\mathcal{I}_n$  (and its size) will be determined based on the problem size n. Each state  $i \in \mathcal{I}_n$  corresponds to an intermediate value dp(i) that needs to be computed, and  $i \prec j$  means that state i needs to be solved before state j. There exists a function  $f_{\mathcal{I}}: \mathcal{I}_n \to \mathcal{I}_n$  to calculate the next state, that is,  $j = f_{\mathcal{I}}(i)$  if j is the next state to solve after i.

• Transition function: The intermediate DP value can be calculated by the transition function f<sub>T</sub> as dp(i) =  $f_T(n, s, \{(j, dp(j)) : j \prec i\})$  where s is the concatenation of the input tokens which corresponds to all elements of all input sequences. Furthermore, this can be formulated as dp(i) =  $f_T(n, \{s_j : j \in D_i\}, \{dp(k) : k \in V_{dp(i)}\})$  where  $D_i$  and  $V_{dp(i)}$  are the sets of input tokens indices and DP values needed to solve state *i* respectively.

• Aggregation function: To produce the final answer, the aggregation function needs to collect the required intermediate DP values and calculate the final result, which can be formalized as anwser =  $f_{\mathcal{A}}(\{dp(i) : i \in \mathcal{A}_n\})$  where  $\mathcal{A}_n$  is the set of DP values needed in the aggregation according to the problem size n.

It should be noted that in the above definition, we use  $s^{(i)}$  to denote the *i*-th input sequence and *s<sub>i</sub>* to denote the *i*-th input token, where *s* is all input tokens transformed from the concatenated input sequences  $(s^{(1)}, s^{(2)}, \ldots, s^{(N)})$ . It can be referenced from Section 4.1 of Feng et al. (2024) for more detailed examples for DP problems. We consider the process by which the Mamba layer defined as Eq (2) gradually generates the solution to DP problems when using CoT. The format of the generated sequence can be written as:

 $s^{(1)} | s^{(2)} | \dots | s^{(N)} | (i_1, dp(i_1)) \dots (i_{|\mathcal{I}_n|}, dp(i_{|\mathcal{I}_n|}))$  final answer

where the input sequence is separated using the symbol | as a delimiter.

**Assumption 2.** Given the input sequences  $s^{(1)}, s^{(2)}, \ldots, s^{(N)}$ , we consider the following constraints for the DP problem:

- For any  $i \in \mathcal{I}_n$ , there exists  $|\mathcal{D}_i| \leq N_s$ ,  $|\mathcal{V}_{dp(i)}| \leq N_{dp}$  and  $|\mathcal{A}(n)| \leq N_{\mathcal{A}}$ .
- The size of the state space  $|\mathcal{I}_n|$ , all elements of input sequences (or equivalently,  $s_i$  for all  $i \in [|s|]$ ), all intermediate DP values  $(dp(i) \text{ for any } i \in \mathcal{I}_n)$ , and the final answer can all be polynomially upper bounded by the problem size n.
- The functions used to solve the DP problem, including the function  $f_{\mathcal{I}}$  to determine the next state, the transition function  $f_{\mathcal{T}}$ , the aggregation function  $f_{\mathcal{A}}$  and  $\mathcal{A}(\mathbf{n})$  can all be approximated with polynomial efficiency by a constant-size MLP (with the SiLU activation function) (Feng et al., 2024).

The first constraint of the Assumption 2 illustrates that the number of input tokens and previous DP values used in the transition function is upper bounded by  $N_s$  and  $N_{dp}$  respectively. In addition, the number of DP values used in aggregation is at most  $N_A$ . The second constraint allows that all involved inputs and outputs used in functions can be represented by the log-precision model. The third constraint allows a constant-sized degenerated Mamba (see Lemma 2) to implement functions required to solve the DP. In fact, due to the first constraint, the sizes of inputs and outputs of these functions will be a constant related to  $\{N_s, N_{dp}, N_A\}$ . Now we present our conclusion as follows:

**Theorem 3** (Perform DP problems with CoT). Considering any DP problem and given input sequences that satisfies Assumption 2, for any integer  $T \in \mathbb{N}$ , there exists several Mamba layers with size O(T), such that the answer generated by the Mamba layers will be correct when the length of the answer is no more than T.

The proof of Theorem 3 can be seen in Appendix A.3. Theorem 3 states that for any DP problem, to generate correct answers with a length no greater than T in CoT, the size of the Mamba layer should scale linearly with the sequence length. Additionally, it can be noted that each step

Models with CoT	Standard TF	Sparse/Linear TF	Mamba
storage	O(1)	$O(\sqrt{T})$	O(T)
inference each step	O(T)	O(T)	O(T)
all steps	$O(T^2)$	$O(T^2)$	$O(T^2)$

Table 1: The comparison of costs for solving any DP problems between Transformers and Mamba.

447 of CoT for Mamba has a complexity of O(T) and thus the total cost is  $O(T^2)$ . In comparison with 448 standard Transformers, Feng et al. (2024) point out that a constant-sized Transformer (O(1) relative 449 to T) can solve any DP problem when equipped with CoT. Thus during each step of CoT inference, the inference cost for Transformer is O(T), making the total cost still  $O(T^2)$ ; While for efficient 450 451 Transformers, especially linear or sparse Transformers, the conclusions obtained by Yang et al. (2024) are also similar: to solve any DP problem, they all require a hidden dimension of  $O(\sqrt{T})$ 452 (and particularly, block size  $B = \Theta(\sqrt{T})$  for sparse transformers), which results in a inference cost 453 of O(T). Consequently, the total cost remains  $O(T^2)$ . Therefore, when equipped with CoT and 454 solving any DP problems, Mamba does not offer additional cost savings compared to Transformers 455 and thus all models are on equal footing in this regard. This comparison can be illustrated in Table 1. 456

457 It seems that this is a frustrating conclusion: similar to efficient Transformers, Mamba also appears 458 not to provide savings in overhead. However, we should note that although there seems to be no 459 shortcut to solving arbitrary DP problems, it may be possible to achieve solutions for slightly sim-460 pler problems using efficient models. We aim to show that just as efficient Transformers can offer advantages when dealing with DP problems with local properties introduced by Yang et al. (2024), 461 Mamba can also similarly address such problems with a smaller overhead compared to standard 462 Transformers. We assume that when solving some DP problem with CoT, the output tokens can be 463 written as  $o_1, o_2, \ldots, o_T$ . If  $o_i = f(\{o_i : i - m \le j \le i\})$  for any  $i \in [T]$ , that is, the calculation 464 for  $o_i$  only depends on at most m preceding intermediate DP values, then we call the DP problem is 465 *m*-locality DP problem. With this assumption, we present the following result: 466

**Theorem 4** (Perform *m*-locality DP problems with CoT). Consider any *m*-locality DP problem and given input sequences that satisfies Assumption 2, for any integer  $T \in \mathbb{N}$ , there exists several Mamba layers with size O(m), such that the answer generated by the Mamba layers will be correct when the length of the answer is no more than T.

The proof of Theorem 4 can be seen in Appendix A.4. Theorem 4 shows that when handling mlocality DP problem with CoT, the needed size of Mamba depends on the problem's locality. When m is much smaller than the total needed length T, the cost for each step becomes a constant O(m); therefore, the total cost becomes O(mT) rather than  $O(T^2)$  leading to savings in cost.

475 476

477

439

6 DISCUSSION

478 In this paper, inspired by the similarity between the SSM module in Mamba and linear attention, 479 we explore Mamba's potential bottlenecks in the COPY operation and show that Mamba with linear 480 size can complete it. Additionally, we present that Mamba has the same cost as standard or efficient 481 Transformers when solving DP problems using CoT. Our findings contribute to a deeper understand-482 ing of Mamba. However, we would like to illustrate that while Mamba may slightly underperform Transformers in certain tasks, it offers advantages in others like sparse parity learning(Park et al., 483 2024) and can achieve comparable performance with lower costs(Gu & Dao, 2023). Therefore, as 484 shown in Park et al. (2024); Waleffe et al. (2024); Wen et al. (2024), exploring hybrid architectures 485 and deeper theoretical analysis for them is a promising direction for future work.

## 486 REFERENCES

524

532

- Ekin Akyürek, Bailin Wang, Yoon Kim, and Jacob Andreas. In-context language learning: Architectures and algorithms. In Ruslan Salakhutdinov, Zico Kolter, Katherine Heller, Adrian Weller, Nuria Oliver, Jonathan Scarlett, and Felix Berkenkamp (eds.), *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pp. 787–812. PMLR, 21–27 Jul 2024.
- Alex Andonian, Quentin Anthony, Stella Biderman, Sid Black, Preetham Gali, Leo Gao, Eric Hallahan, Josh Levy-Kramer, Connor Leahy, Lucas Nestler, Kip Parker, Michael Pieler, Jason Phang, Shivanshu Purohit, Hailey Schoelkopf, Dashiell Stander, Tri Songz, Curt Tigges, Benjamin Thérien, Phil Wang, and Samuel Weinbach. GPT-NeoX: Large Scale Autoregressive Language Modeling in PyTorch, 9 2023. URL https://www.github.com/eleutherai/ gpt-neox.
- Simran Arora, Sabri Eyuboglu, Aman Timalsina, Isys Johnson, Michael Poli, James Zou, Atri Rudra, and Christopher Ré. Zoology: Measuring and improving recall in efficient language models. *arXiv preprint arXiv:2312.04927*, 2023.
- Iz Beltagy, Matthew E Peters, and Arman Cohan. Longformer: The long-document transformer.
   *arXiv preprint arXiv:2004.05150*, 2020.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,
  Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are
  few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- Rewon Child, Scott Gray, Alec Radford, and Ilya Sutskever. Generating long sequences with sparse
   transformers. *arXiv preprint arXiv:1904.10509*, 2019.
- Krzysztof Choromanski, Valerii Likhosherstov, David Dohan, Xingyou Song, Andreea Gane, Tamas
  Sarlos, Peter Hawkins, Jared Davis, Afroz Mohiuddin, Lukasz Kaiser, et al. Rethinking attention
  with performers. *arXiv preprint arXiv:2009.14794*, 2020.
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm:
  Scaling language modeling with pathways. *Journal of Machine Learning Research*, 24(240):
  1–113, 2023.
- Tri Dao. Flashattention-2: Faster attention with better parallelism and work partitioning. *arXiv* preprint arXiv:2307.08691, 2023.
- Tri Dao and Albert Gu. Transformers are ssms: Generalized models and efficient algorithms through
   structured state space duality. *arXiv preprint arXiv:2405.21060*, 2024.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas
   Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An
   image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- Guhao Feng, Bohang Zhang, Yuntian Gu, Haotian Ye, Di He, and Liwei Wang. Towards revealing
   the mystery behind chain of thought: a theoretical perspective. *Advances in Neural Information Processing Systems*, 36, 2024.
- Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang, Horace He, Anish Thite, Noa Nabeshima, et al. The pile: An 800gb dataset of diverse text for language modeling. *arXiv preprint arXiv:2101.00027*, 2020.
- Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. *arXiv preprint arXiv:2312.00752*, 2023.
- 539 Albert Gu, Karan Goel, and Christopher Ré. Efficiently modeling long sequences with structured state spaces. *arXiv preprint arXiv:2111.00396*, 2021.

540 Albert Gu, Karan Goel, Ankit Gupta, and Christopher Ré. On the parameterization and initialization 541 of diagonal state space models. Advances in Neural Information Processing Systems, 35:35971-542 35983, 2022. 543 Ankit Gupta, Albert Gu, and Jonathan Berant. Diagonal state spaces are as effective as structured 544 state spaces. Advances in Neural Information Processing Systems, 35:22982–22994, 2022. 546 Dongchen Han, Ziyi Wang, Zhuofan Xia, Yizeng Han, Yifan Pu, Chunjiang Ge, Jun Song, Shiji 547 Song, Bo Zheng, and Gao Huang. Demystify mamba in vision: A linear attention perspective. arXiv preprint arXiv:2405.16605, 2024. 548 549 Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and 550 Jacob Steinhardt. Measuring massive multitask language understanding. arXiv preprint 551 arXiv:2009.03300, 2020. 552 Samy Jelassi, David Brandfonbrener, Sham M Kakade, and Eran Malach. Repeat after me: Trans-553 formers are better than state space models at copying. arXiv preprint arXiv:2402.01032, 2024. 554 555 Angelos Katharopoulos, Apoorv Vyas, Nikolaos Pappas, and François Fleuret. Transformers are 556 rnns: Fast autoregressive transformers with linear attention. pp. 5156–5165, 2020. Jacob Devlin Ming-Wei Chang Kenton and Lee Kristina Toutanova. Bert: Pre-training of deep 558 bidirectional transformers for language understanding. In *Proceedings of naacL-HLT*, volume 1, 559 pp. 2, 2019. 560 561 Zhiyuan Li, Hong Liu, Denny Zhou, and Tengyu Ma. Chain of thought empowers transformers to 562 solve inherently serial problems. arXiv preprint arXiv:2402.12875, 2024. 563 I Loshchilov. Decoupled weight decay regularization. arXiv preprint arXiv:1711.05101, 2017. 564 565 William Merrill and Ashish Sabharwal. The expressive power of transformers with chain of 566 thought. arXiv preprint arXiv:2310.07923, 2023a. 567 William Merrill and Ashish Sabharwal. The parallelism tradeoff: Limitations of log-precision trans-568 formers. Transactions of the Association for Computational Linguistics, 11:531–545, 2023b. 569 570 William Merrill, Jackson Petty, and Ashish Sabharwal. The illusion of state in state-space models. 571 arXiv preprint arXiv:2404.08819, 2024. 572 Erxue Min, Runfa Chen, Yatao Bian, Tingyang Xu, Kangfei Zhao, Wenbing Huang, Peilin Zhao, 573 Junzhou Huang, Sophia Ananiadou, and Yu Rong. Transformer for graphs: An overview from 574 architecture perspective. arXiv preprint arXiv:2202.08455, 2022. 575 Jongho Park, Jaeseung Park, Zheyang Xiong, Nayoung Lee, Jaewoong Cho, Samet Oymak, Kang-576 wook Lee, and Dimitris Papailiopoulos. Can mamba learn how to learn? A comparative study 577 on in-context learning tasks. In Proceedings of the 41st International Conference on Machine 578 Learning, volume 235 of Proceedings of Machine Learning Research, pp. 39793–39812. PMLR, 579 21-27 Jul 2024. 580 581 Bo Peng, Eric Alcaide, Quentin Anthony, Alon Albalak, Samuel Arcadinho, Stella Biderman, 582 Huanqi Cao, Xin Cheng, Michael Chung, Matteo Grella, et al. Rwkv: Reinventing rnns for 583 the transformer era. arXiv preprint arXiv:2305.13048, 2023. 584 Noam Shazeer. Glu variants improve transformer. arXiv preprint arXiv:2002.05202, 2020. 585 586 Jerome Sieber, Carmen Amo Alonso, Alexandre Didier, Melanie N Zeilinger, and Antonio Orvieto. Understanding the differences in foundation models: Attention, state space models, and recurrent neural networks. arXiv preprint arXiv:2405.15731, 2024. 588 589 Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan, Wen Bo, and Yunfeng Liu. Roformer: En-590 hanced transformer with rotary position embedding. Neurocomputing, 568:127063, 2024. Yutao Sun, Li Dong, Shaohan Huang, Shuming Ma, Yuqing Xia, Jilong Xue, Jianyong Wang, and 592 Furu Wei. Retentive network: A successor to transformer for large language models. arXiv 593

preprint arXiv:2307.08621, 2023.

- Yi Tay, Mostafa Dehghani, Dara Bahri, and Donald Metzler. Efficient transformers: A survey. ACM *Comput. Surv.*, 55(6):109:1–109:28, 2023.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée
  Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and
  efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
   Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. Advances in neural information processing systems, 30, 2017.
  - Roger Waleffe, Wonmin Byeon, Duncan Riach, Brandon Norick, Vijay Korthikanti, Tri Dao, Albert Gu, Ali Hatamizadeh, Sudhakar Singh, Deepak Narayanan, et al. An empirical study of mambabased language models. *arXiv preprint arXiv:2406.07887*, 2024.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny
   Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837, 2022.
- Kaiyue Wen, Xingyu Dang, and Kaifeng Lyu. Rnns are not transformers (yet): The key bottleneck on in-context retrieval. *arXiv preprint arXiv:2402.18510*, 2024.
- Rui Xu, Shu Yang, Yihui Wang, Bo Du, and Hao Chen. A survey on vision mamba: Models, applications and challenges. *arXiv preprint arXiv:2404.18861*, 2024.
- Kai Yang, Jan Ackermann, Zhenyu He, Guhao Feng, Bohang Zhang, Yunzhen Feng, Qiwei Ye, Di He, and Liwei Wang. Do efficient transformers really save computation? *arXiv preprint arXiv:2402.13934*, 2024.
- Lianghui Zhu, Bencheng Liao, Qian Zhang, Xinlong Wang, Wenyu Liu, and Xinggang Wang. Vision mamba: Efficient visual representation learning with bidirectional state space model. *arXiv* preprint arXiv:2401.09417, 2024.

### A APPENDIX

### A.1 PROOF OF THEOREM 1

**Theorem 5** (Perform COPY operation). Given a input sequence  $x_1, x_2, ..., x_N$  and for any  $\epsilon > 0$ , there exists an SSM module with constant size that can approximate the COPY operation defined above when the following condition is satisfied for  $i \in [N]$ :

$$\rho \ge \left(1 - \frac{\epsilon}{2M \|\boldsymbol{\Delta}\|_{\infty}}\right) \frac{1}{\alpha_{pos(i)}} + \delta \left(\frac{a_{\max,<}}{1 - a_{\max,<}} + \frac{\left(\frac{1}{a_{\min,>}}\right)^{L-1} - 1}{1 - a_{\min,>}}\right),\tag{8}$$

where  $a_{\max,<} = \max_{1 \le j < pos(i)} a_j$  and  $a_{\min,>} = \min_{pos(i) < j \le i} a_j$ . Moreover, in such cases, there exists  $\delta \le \frac{\epsilon}{2(L-1)M\|\mathbf{\Delta}\|_{\infty}}$  when  $L \ge 2$ .

*Proof.* We firstly show that given the *i*-th input  $x_i$ , a SSM module can retrieve the most relevant historical record  $v_{pos(i)}$  from the hidden state and perform the defined COPY operation under the condition illustrated in our theorem. To achieve this, recalling that  $v_j = \Delta_j \odot x_j$  in Eq (6), we have that

$$\left\|\boldsymbol{y}_{i}-\boldsymbol{v}_{pos(i)}\right\|_{\infty}=\left\|\sum_{j=1}^{i}\alpha_{j}(\boldsymbol{\Delta}_{j}\odot\boldsymbol{x}_{j})\boldsymbol{b}_{j}^{T}\boldsymbol{c}_{i}-\boldsymbol{v}_{pos(i)}\right\|_{\infty}=\left\|\sum_{j=1}^{i}\alpha_{j}(\boldsymbol{c}_{i}^{T}\boldsymbol{b}_{j})\boldsymbol{v}_{j}-\boldsymbol{v}_{pos(i)}\right\|_{\infty}$$
(9)

$$= \left\| \sum_{j \neq pos(i)} \alpha_j (\boldsymbol{c}_i^T \boldsymbol{b}_j) \boldsymbol{v}_j + \alpha_{pos(i)} (\boldsymbol{c}_i^T \boldsymbol{b}_{pos(i)}) \boldsymbol{v}_{pos(i)} - \boldsymbol{v}_{pos(i)} \right\|_{\infty}$$
(10)

 $= \left\| \sum_{j \neq pos(i)} \alpha_j(\boldsymbol{c}_i^T \boldsymbol{b}_j) \boldsymbol{v}_j + \left[ \alpha_{pos(i)}(\boldsymbol{c}_i^T \boldsymbol{b}_{pos(i)}) - 1 \right] \boldsymbol{v}_{pos(i)} \right\|_{\infty}$ (11)

$$\leq M \|\boldsymbol{\Delta}\|_{\infty} \left( \sum_{j \neq pos(i)} \alpha_j |\boldsymbol{c}_i^T \boldsymbol{b}_j| + 1 - \alpha_{pos(i)} \boldsymbol{c}_i^T \boldsymbol{b}_{pos(i)} \right), \tag{12}$$

where in (12) we use the fact that  $\|\boldsymbol{x}_i\|_{\infty} \leq M$  and  $\alpha_{pos(i)}|\boldsymbol{c}_i^T \boldsymbol{b}_{pos(i)}| \leq 1$  (from the second condition in Assumption 1). Thus, to prove  $\|\boldsymbol{y}_i - \boldsymbol{v}_{pos(i)}\|_{\infty} \leq \epsilon$ , we can show that

$$\boldsymbol{c}_{i}^{T}\boldsymbol{b}_{pos(i)} \geq \left(1 - \frac{\epsilon}{M\|\boldsymbol{\Delta}\|_{\infty}}\right) \frac{1}{\alpha_{pos(i)}} + \sum_{j \neq pos(i)} \frac{\alpha_{j}}{\alpha_{pos(i)}} |\boldsymbol{c}_{i}^{T}\boldsymbol{b}_{j}|.$$
(13)

Recalling that  $\alpha_j = \prod_{k=j+1}^i a_k$ , we have

$$\frac{\alpha_j}{\alpha_{pos(i)}} = \begin{cases} \prod_{k=j+1}^{pos(i)} a_k = a_{pos(i)} a_{pos(i)-1} \dots a_{j+1}, & \text{when } j < pos(i) \\ \prod_{k=pos(i)+1}^j \frac{1}{a_k} = \frac{1}{a_j a_{j-1} \dots a_{pos(i)+1}}, & \text{when } j > pos(i). \end{cases}$$
(14)

Then we can consider the second term on the right side of Inequality (13) as

$$\sum_{j \neq pos(i)} \frac{\alpha_j}{\alpha_{pos(i)}} |\boldsymbol{c}_i^T \boldsymbol{b}_j| = \sum_{j \notin S_i} \frac{\alpha_j}{\alpha_{pos(i)}} |\boldsymbol{c}_i^T \boldsymbol{b}_j| + \sum_{j \in S_i, j \neq pos(i)} \frac{\alpha_j}{\alpha_{pos(i)}} |\boldsymbol{c}_i^T \boldsymbol{b}_j|$$
(15)

For the first term on the right side, we have 

$$\sum_{j \notin S_i} \frac{\alpha_j}{\alpha_{pos(i)}} \left| \boldsymbol{c}_i^T \boldsymbol{b}_j \right| \le \delta \left( \sum_{j \notin S_i, j < pos(i)} \frac{\alpha_j}{\alpha_{pos(i)}} + \sum_{j \notin S_i, j > pos(i)} \frac{\alpha_j}{\alpha_{pos(i)}} \right)$$
(16)

$$\leq \delta \left( \sum_{j=1}^{pos(i)-1} \frac{\alpha_j}{\alpha_{pos(i)}} + \sum_{j=pos(i)+1}^{i} \frac{\alpha_j}{\alpha_{pos(i)}} \right)$$
(17)

710  
711  
712  
713
$$\leq \delta \left( \sum_{k=1}^{pos(i)-1} (a_{\max,<})^k + \sum_{k=1}^{i-pos(i)} \left( \frac{1}{a_{\min,>}} \right)^k \right)$$
(18)

714  
715  
716  
717  
717  

$$\leq \delta \left( \frac{a_{\max,<}(1-a_{\max,<}^{pos(i)-1})}{1-a_{\max,<}} + \frac{\left(\frac{1}{a_{\min,>}}\right)^{i-pos(i)}-1}{1-a_{\min,>}} \right)$$
(19)

where  $a_{\max,<} = \max_{1 \le j < pos(i)} a_j$  and  $a_{\min,>} = \min_{pos(i) < j \le i} a_j$ . In (16) we use the assumption that  $|c_i^T b_j| le\delta$  for  $j \notin S_i$ ; in (17) we use  $\frac{\alpha_j}{\alpha_{pos(i)}} > 0$  for all  $j \le i$ ; in (18), we use the fact that  $\frac{\alpha_j}{\alpha_{pos(i)}} \leq (a_{\max,<})^{pos(i)-j} \text{ for } j < pos(i) \text{ and } \frac{\alpha_j}{\alpha_{pos(i)}} \leq \left(\frac{1}{a_{\min,>}}\right)^{j-pos(i)} \text{ for } j > pos(i); \text{ in (19), we use the formula for the sum of a geometric series.}$ 

Furthermore, considering that the vector  $v_{pos(i)}$  to be copied must exist in the L-local matching set  $S_i$  so there is  $i - L + 1 \le pos(i) \le i$ , we have the following 

$$\sum_{j \notin S_i} \frac{\alpha_j}{\alpha_{pos(i)}} \left| \boldsymbol{c}_i^T \boldsymbol{b}_j \right| \le \delta \left( \frac{a_{\max,<} (1 - a_{\max,<}^{pos(i)-1})}{1 - a_{\max,<}} + \frac{\left(\frac{1}{a_{\min,>}}\right)^{L-1} - 1}{1 - a_{\min,>}} \right)$$
(20)

> $\leq \delta \left( \frac{a_{\max,<}}{1 - a_{\max,<}} + \frac{\left(\frac{1}{a_{\min,>}}\right)^{L-1} - 1}{1 - a_{\min,>}} \right).$ (21)

In (20), we use the fact that  $a_{\max,<} \in (0,1)$  while  $\frac{1}{a_{\min,>}} > 1$ ; in (21) we just ignore the term  $a_{\max,<}^{pos(i)-1}$  for simplicity (in fact, we can find that when *i* is sufficiently large, the effect of this term can be neglected). 

Meanwhile, with Assumption 1, we can consider the second term on the right side of Inequality (15) as 

$$\sum_{\substack{\in S_i, j \neq pos(i) \\ \alpha_{pos(i)}}} \frac{\alpha_j}{\alpha_{pos(i)}} \left| \boldsymbol{c}_i^T \boldsymbol{b}_j \right| \le \frac{1}{\alpha_{pos(i)}} - \boldsymbol{c}_i^T \boldsymbol{b}_{pos(i)}.$$
(22)

Thus, considering (13), (21) and (22), to show  $\|y_i - o_i\|_{\infty} \leq \epsilon$ , the following condition should be satisfied 

$$\boldsymbol{c}_{i}^{T}\boldsymbol{b}_{pos(i)} \geq \left(1 - \frac{\epsilon}{M\|\boldsymbol{\Delta}\|_{\infty}}\right) \frac{1}{\alpha_{pos(i)}} + \delta\left(\frac{a_{\max,<}}{1 - a_{\max,<}} + \frac{\left(\frac{1}{a_{\min,>}}\right)^{L-1} - 1}{1 - a_{\min,>}}\right)$$
(23)

$$+\frac{1}{\alpha_{pos(i)}} - \boldsymbol{c}_i^T \boldsymbol{b}_{pos(i)}.$$
(24)

After reformulating, there exists

j

$$\boldsymbol{c}_{i}^{T}\boldsymbol{b}_{pos(i)} \geq \left(1 - \frac{\epsilon}{2M\|\boldsymbol{\Delta}\|_{\infty}}\right) \frac{1}{\alpha_{pos(i)}} + \frac{\delta}{2} \left(\frac{a_{\max,<}}{1 - a_{\max,<}} + \frac{\left(\frac{1}{a_{\min,>}}\right)^{L-1} - 1}{1 - a_{\min,>}}\right).$$
(25)

756 757 758 Recalling that  $c_i b_{pos(i)}^T \ge \rho$ , we can find that if we set the lower bound of  $\rho$  as the right side of Inequality (25), we will have  $\|\boldsymbol{y}_i - \boldsymbol{o}_i\| \le \epsilon$ .

Moreover, in such case, with the second condition of Assumption 1, we have  $c_i^T b_{pos(i)} \leq \frac{1}{\alpha_{pos(i)}}$ thus there is

763 764 765

$$\frac{1}{\alpha_{pos(i)}} \ge \left(1 - \frac{\epsilon}{2M \|\mathbf{\Delta}\|_{\infty}}\right) \frac{1}{\alpha_{pos(i)}} + \frac{\delta}{2} \left(\frac{a_{\max,<}}{1 - a_{\max,<}} + \frac{\left(\frac{1}{a_{\min,>}}\right)^{L-1} - 1}{1 - a_{\min,>}}\right).$$
(26)

By reformulating, when  $L \ge 2$ , we can get that

$$\delta \le \frac{\epsilon}{M \|\mathbf{\Delta}\|_{\infty}} \left( \frac{a_{\max,<}}{1 - a_{\max,<}} + \frac{\left(\frac{1}{a_{\min,>}}\right)^{L-1} - 1}{1 - a_{\min,>}} \right)^{-1}$$
(27)

770 771 772

773

774 775

776

777

$$\leq \frac{\epsilon}{M \|\boldsymbol{\Delta}\|_{\infty}} \left( \frac{\left(\frac{1}{a_{\min,>}}\right)^{L-1} - 1}{1 - a_{\min,>}} \right)^{-1}$$
(28)

$$\leq \frac{\epsilon}{M \|\boldsymbol{\Delta}\|_{\infty}} \cdot \frac{1}{L-1},\tag{29}$$

where in (29) we use the inequality  $\frac{1-a}{(\frac{1}{a})^x-1} \leq \frac{1}{x}$  for 0 < a < 1 and  $x \geq 1$  where a and x are replaced by  $a_{min,>}$  and L-1 in (28) respectively when  $L \geq 2$ .

We have shown above that an SSM module can perform the defined COPY operation under the conditions described in Theorem 1. To prove that a Mamba block can do the same, we only need to demonstrate that a Mamba block defined by Eq (3) degenerates into an SSM module. In fact, we only need to deactivate the gating branch to achieve this. For example, we can set  $W_1 = W_3 =$  $I, b_1 = 0, W_2 = O$  and  $b_2 = k1$  where the constant k satisfies  $\sigma(k) = 1$ . Thus, we complete our proof.

787 788 A.2 PROOF OF THEOREM 2

**Theorem 6** (Perform COPY operation with linear-scaling size). Given sequence  $x_1, x_2, ..., x_N \in [-M, M]^d$ , there exists a Mamba block with size O(N) that can perform the defined COPY operation, that is,  $y_i = o_i$  for any  $i \in [N]$ . Moreover, the  $l_{\infty}$  norm of the Mamba block parameters is upper bounded by O(ploy(M, N)).

793

**Proof.** Recalling that the output of SSM module in Mamba can be rewritten in the form of Eq (6), where  $(\Delta_j \odot x_j)$ ,  $b_j$ ,  $c_i$  corresponds to  $v_j$ ,  $k_j$  and  $q_i$  respectively. Our intuition is to store all the information of  $v_i$  from our history in the hidden state space of size O(N) (similar to the KV cache in attention format), and then use the appropriate  $c_i$  as the query for retrieval. We can set A = O so that Eq (6) further transforms in a way that does not forget historical information, that is,  $y_i = \sum_{j=1}^{i} (\Delta_j \odot x_j) b_j^T c_i = \sum_{j=1}^{i} v_j b_j^T c_i$ .

Let  $\tilde{x}_i = [x_i, e_i, e_{pos(i)}] \in \mathbb{R}^{d+2N}$  where  $e_i \in \mathbb{R}^N$  denote the one-hot vector where only the *i*-th 801 value is 1. We use  $e_i$ ,  $e_{pos(i)}$  to denote the current position and the position of historical token 802 we want to copy respectively. Then, we construct  $W_b = [O_{N \times d}, I_N, O_N] \in \mathbb{R}^{N \times (d+2N)}$  so that 803  $ilde{b}_i = W_b ilde{x}_i = e_i$ . Then, at the *i*-th step, the information newly recorded in the state space will be 804  $v_i \tilde{b}_i^T = v_i e_i^T \in \mathbb{R}^{d \times N}$  and the updated state space will be  $H_i = H_{i-1} + \sum_{j=1}^{i-1} v_j \tilde{b}_j^{\hat{T}} + v_i \tilde{b}_i^T =$ 805 806  $[v_1, v_2, \dots, v_i, O_{d \times (N-i)}]$  thus at the last step, we can record all historical information in the state 807 space by  $\boldsymbol{H}_T = \sum_{j=1}^N \boldsymbol{v}_j \tilde{\boldsymbol{b}}_j^T = [\boldsymbol{v}_1, \boldsymbol{v}_2, \dots, \boldsymbol{v}_N]$ . Then, at the output process, we can construct  $\boldsymbol{W}_{\boldsymbol{c}} = [\boldsymbol{O}_{N \times d}, \boldsymbol{O}_N, \boldsymbol{I}_N] \in \mathbb{R}^{N \times (d+2N)}$  so that  $\tilde{\boldsymbol{c}}_i = \boldsymbol{W}_{\boldsymbol{c}} \tilde{\boldsymbol{x}}_i = \boldsymbol{e}_{pos(i)}$ . Thus, the output will be 808 809  $oldsymbol{y}_i = oldsymbol{H}_i oldsymbol{ ilde{c}}_i = \sum_{i=1}^i oldsymbol{v}_j oldsymbol{b}_j^T oldsymbol{e}_{pos(i)} = oldsymbol{v}_{pos(i)}.$ 

At the same time, we note that in the above process, the vectors  $e_i, e_{pos(i)}$  to denote position in  $\tilde{x}_i$ are sparse. In fact, we only need to use two indices  $p_i = i$  and  $p_{pos(i)} = pos(i)$  to store them thus the total size to store all indices is O(N). Additionally,  $W_b, W_c$  are also sparse so we require at most O(N) space to store these two matrices. Therefore, the model size we need is O(Nd) that scales linearly with the length N. Similar to the proof of Theorem 1, we can degenerate the Mamba block into the aforementioned SSM module by deactivating the gating branch. In addition, we note that  $x_i \in [-M, M]^d$  and  $p_i, p_{pos(i)} \in [1, N]$  thus the largest value involved in the aforemen-tioned process will not exceed  $NM^2$  (the largest value in hidden states), which is upper bounded by O(ploy(M, N)). Thus, we complete our proof. 

### A.3 PROOF OF THEOREM 3

In this part, we first present the necessary lemmas before completing the proof of Theorem 3. In
fact, these lemmas are very similar to those presented by Feng et al. (2024) in Appendix C.1 regarding MLP. The main difference is that we need to degenerate Mamba blocks to MLP and consider
different activation functions (SiLU for Mamba instead of GELU). Thus, we only provide detailed
proofs of these relevant lemmas when necessary.

**Lemma 1** (Perform multiplication). For any  $\epsilon > 0$  and M > 0, there exists Mamba block parameters with  $l_{\infty}$  norm upper bounded by  $O(poly(M, 1/\epsilon))$  such that  $|f(a, b) - ab| \le \epsilon$  holds for all  $a, b \in [-M, M]$ .

*Proof.* We first show that a two-layer MLP using the SiLU activation function can achieve the above operation. We use the same construction as in Lemma C.1. in Feng et al. (2024), except that we use the SiLU activation function instead of GELU. Specifically, let  $g : \mathbb{R}^2 \to \mathbb{R}$  be a two-layer MLP with SiLU activation, and the hidden dimension is 4, then we can construct f as

$$g(a,b) = \frac{\lambda^2}{2} \left( \sigma\left(\frac{a+b}{\lambda}\right) + \sigma\left(\frac{-a-b}{\lambda}\right) - \sigma\left(\frac{a-b}{\lambda}\right) - \sigma\left(\frac{-a+b}{\lambda}\right) \right), \tag{30}$$

where  $\lambda$  is a scaling factor. In addition, considering  $\sigma(x) = \frac{x}{1+e^{-x}}$ ,  $\sigma'(x) = \frac{1+(x+1)e^{-x}}{(1+e^{-x})^2}$ ,  $\sigma''(x) = \frac{e^{-x}(2+2e^{-x}+xe^{-x}-x)}{(1+e^{-x})^3}$ , we have  $\sigma(0) = 0$ ,  $\sigma'(0) = \frac{1}{2}$ ,  $\sigma''(0) = \frac{1}{2}$ . Then, using the Taylor expansion with the Lagrange remainder, we can obtain that

 $\sigma\left(\frac{a+b}{\lambda}\right) + \sigma\left(\frac{-a-b}{\lambda}\right) - \sigma\left(\frac{a-b}{\lambda}\right) - \sigma\left(\frac{-a+b}{\lambda}\right)$ 

 $4 (2M)^2$ 

where  $R_2$  is the second-order remainder term. Assuming that  $\lambda > 2M$ , we have  $\left|\frac{\pm a \pm b}{\lambda}\right| < \frac{2M}{\lambda} < 1$ and then

 $=\frac{1}{2!}\frac{1}{2}\left(\left(\frac{a+b}{\lambda}\right)^2+\left(\frac{-a-b}{\lambda}\right)^2-\left(\frac{a-b}{\lambda}\right)^2-\left(\frac{-a+b}{\lambda}\right)^2\right)+R_2=\frac{2ab}{\lambda^2}+R_2,$ 

$$|R_2| \leq \frac{1}{3!} \left(\frac{1}{\lambda}\right) \max_{x \in [-1,1]} |\sigma^{m}(x)|$$

$$= \frac{4}{3!} \left(\frac{2M}{\lambda}\right)^2 \max_{x \in [-1,1]} \left|\frac{(x-3)e^{-x} - 4xe^{-2} + (x+3)e^{-3x}}{(1+e^{-x})^4}\right|$$

$$\leq \frac{4}{3!} \left(\frac{2M}{\lambda}\right)^2 \frac{4e + 4e^2 + 4e^3}{(1+e^{-1})^4}$$

$$\leq \frac{4}{3!} \frac{8M^3}{\lambda^3} \frac{81}{2}$$

$$= \frac{216M^3}{\lambda^3}.$$

Thus if we set  $\lambda \geq \frac{216M^3}{2\epsilon}$  we will have  $|g(a,b) - ab| \leq \frac{\lambda^2}{2} |R_2| \leq \epsilon$ .

Then, we note that a Mamba block f defined as Eq (3) can degenerate into the above MLP g by deactivating its SSM branch. Specifically, we only need to set  $W_1$  to be zeros and b = 1 so that the

input of the SSM branch is a constant 1, that is,

866

872

887

895 896 897

899

909

$$\boldsymbol{f}(\boldsymbol{x}) = \boldsymbol{W}_3 \cdot \mathrm{SSM}(1) \odot \sigma(\boldsymbol{W}_2 \boldsymbol{x} + \boldsymbol{b}_2).$$

In the SSM module, we can set  $W_b$  to be zeros, that is, no new information will be retained in the hidden state. Following this, we set d = 1 resulting that given x = 1, we have  $SSM(1) = y = c^T H + d \odot x = c^T 1 + 1 \odot 1 = 1$ . Thus the SSM branch can be deactivated and the Mamba block will degenerate into a two layer MLPs, that is,

$$\boldsymbol{f}(\boldsymbol{x}) = \boldsymbol{W}_3 \sigma(\boldsymbol{W}_2 \boldsymbol{x} + \boldsymbol{b}_2). \tag{31}$$

Furthermore, given  $\boldsymbol{x} = [a, b]$ , we can set  $\boldsymbol{W}_2 \in \mathbb{R}^{4 \times 2}$  and  $\boldsymbol{W}_3 \in \mathbb{R}^{4 \times 1}$  to meet the two-layer MLP g as Eq (30). Additionally, we note that all parameters of this Mamba block can be upper bounded by  $O(poly(M, 1/\epsilon))$  under the  $l_{\infty}$  norm. Thus, we complete our proof.

**Remark 1.** It should be noted that here we have only provided one possible construction and this is not unique. For example, in the process of deactivating the SSM branch, we could also choose to make  $\Delta$  sufficiently large and correspondingly  $\tilde{A}$  sufficiently small with  $A \leq 0$  so that the hidden states approximates zeros. In fact, the expressive power of an Mamba block with two branches should be stronger than that of a two-layer MLP since it already encompasses the latter. Nevertheless, we still provide one possible construction here.

**Lemma 2** (Approximate two-layer MLPs with ReLU). Let  $g : \mathbb{R}^{d_1} \to \mathbb{R}^{d_2}$  be a two-layer MLP with ReLU activation, and all parameters are upper bounded by M. Then, for any  $\epsilon > 0$ , there exists a Mamba block f and parameters upper bounded by  $O(poly(M, 1/\epsilon))$  in the  $l_{\infty}$  norm, such that for all  $x \in \mathbb{R}^{d_1}$ , we have  $||f(x) - g(x)||_{\infty} \le \epsilon$ .

**Proof.** Similar to Lemma 1, once again, we deactivate the SSM branch, causing a Mamba block to degenerate into the form of Eq 31. Considering a two-layer MLP with a ReLU activation function denoted as  $g(x) = \overline{W}_3 \text{ReLU}(\overline{W}_2 x)$  where  $\overline{W}_2 \in \mathbb{R}^{d \times d_1}$  and  $\overline{W}_3 \in \mathbb{R}^{d_2 \times d}$ , we can set similar parameters for the degenerated Mamba block, that is, we consider  $W_2 = \lambda \overline{W}_2$ ,  $W_3 = \frac{1}{\lambda} \overline{W}_3$ in Eq (31) where  $\lambda$  is some large constant. In order to prove the lemma, we need to show that  $\|f(x) - g(x)\|_{\infty} \le \epsilon$  with some  $\lambda$  upper bounded by  $O(ploy(M, 1/\epsilon))$ .

Considering a scalar  $z \in \mathbb{R}$ , we firstly consider the upper bound of the following equation:

$$\left| \operatorname{ReLU}(z) - \frac{1}{\lambda} \operatorname{SiLU}(\lambda z) \right| = \left| \max(z, 0) - \frac{z}{1 + e^{-\lambda z}} \right| = \frac{|z|}{e^{\lambda |z|} + 1} \le \frac{1}{\lambda},$$

where we use the fact that  $e^x + 1 > x$  for any  $x \ge 0$ . Then, let  $z = \overline{W}_2 x$ , we can show that for any  $z \in \mathbb{R}^d$ ,

$$\left\|\overline{\boldsymbol{W}}_{3}\operatorname{ReLU}(\boldsymbol{z}) - \frac{1}{\lambda}\overline{\boldsymbol{W}}_{3}\operatorname{SiLU}(\lambda\boldsymbol{z})\right\|_{\infty} \leq \left\|\overline{\boldsymbol{W}}_{3}\right\|_{\infty} \left\|\operatorname{ReLU}(\boldsymbol{z}) - \frac{1}{\lambda}\operatorname{SiLU}(\lambda\boldsymbol{z})\right\|_{\infty}$$
(32)

$$\leq Md \left\| \operatorname{ReLU}(\boldsymbol{z}) - \frac{1}{\lambda} \operatorname{SiLU}(\lambda \boldsymbol{z}) \right\|_{\infty}$$
 (33)

$$\leq Md \max_{z \in \mathbb{R}} \left| \operatorname{ReLU}(z) - \frac{1}{\lambda} \operatorname{SiLU}(\lambda z) \right|$$
 (34)

$$\leq \frac{Md}{\lambda}.$$
(35)

Then, if we set  $\lambda > \frac{Md}{\epsilon}$ , we will have  $\|f(x) - g(x)\|_{\infty} \le \epsilon$  and all parameters of the Mamba block is upper bounded by  $O(ploy(M, 1/\epsilon))$ . Thus, we complete our proof.

913 Remark 2. We have proven that a Mamba block can approximate a two-layer MLP with ReLU
914 activation function, and since the latter can perform many basic operations, including linear trans915 formations and selection operations as constructed in Lemma C.3 and Lemma C.5 in Feng et al.
916 (2024), we can use Lemma 2 to adopt the same construction, enabling the Mamba block to perform
917 these operations. We present the following colloary more specifically, and the detailed proof can be
found in the above mentioned part in Feng et al. (2024).

918 Lemma 3 (Perform linear transformation, easily derived from Lemma 2 and Lemma C.3 in Feng 919 et al. (2024)). Let  $W \in \mathbb{R}^{d_2 \times d_1}$  be any matrix used for implementing linear transformations upper 920 bounded by M and  $f: \mathbb{R}^{d_1} \to \mathbb{R}^{d_2}$  be a Mamba block. Then, for any  $\epsilon > 0$ , there exist Mamba 921 block parameters with  $l_{\infty}$  norm bounded by  $O(poly(M, 1/\epsilon))$ , such that for any  $x \in \mathbb{R}^{d_1}$ , we have 922  $\|\boldsymbol{f}(\boldsymbol{x}) - \boldsymbol{W}\boldsymbol{x}\|_{\infty} \leq \epsilon.$ 

923 Lemma 4 (Perform select operation, easily derived from Lemma 2 and Lemma C.4 in Feng et al. 924 (2024)). Define the selection function  $g : \mathbb{R}^d \times \mathbb{R}^d \times \mathbb{R} \to \mathbb{R}^d$  as follows: 925

$$g(\boldsymbol{x}, \boldsymbol{y}, t) = \begin{cases} \boldsymbol{x} & if \ t > 0\\ \boldsymbol{y} & if \ t < 0 \end{cases}$$
(36)

928 Let  $f: \mathbb{R}^d \times \mathbb{R}^d \times \mathbb{R} \to \mathbb{R}^d$  be a Mamba block. Then, for any  $\epsilon > 0$ ,  $\alpha > 0$ , and M > 0, there 929 exist Mamba parameters with  $l_{\infty}$  norm bounded by  $O(poly(M, 1/\alpha, 1/\epsilon))$ , such that for all  $x \in$ 930  $[-M, M]^d$ ,  $\boldsymbol{y} \in [-M, M]^d$ , and  $t \in [-\infty, -\alpha] \cup [\alpha, +\infty]$ , we have  $\|\boldsymbol{f}(\boldsymbol{x}, \boldsymbol{y}, t) - g(\boldsymbol{x}, \boldsymbol{y}, t)\|_{\infty} \leq \epsilon$ .

932 Next, we show that one Mamba layer or several Mamba layers can implement indicator functions 933 through the select operation. We mainly focus on the usual indicator functions  $\mathbb{I}[a \neq b], \mathbb{I}[a > b]$ and  $\mathbb{I}[a < b]$ . 934

935 **Lemma 5** (Perform indicator function). Define the indicator function  $\mathbb{I}(a, b, \circ) : \mathbb{R}^2 \times \{\neq, >, <$ 936  $\} \rightarrow \{0,1\}$  where  $a, b \in [-M, M]$ . The output of the function will be 1 if  $a \circ b$  is satisfied otherwise 937 the output will be 0. Let  $f : \mathbb{R}^2 \to \mathbb{R}$  be a Mamba block. Then, for any  $\epsilon > 0$ , there exist Mamba 938 parameters with  $l_{\infty}$  norm upper bounded by  $O(poly(M, 1/\epsilon))$ , such that for any  $a, b \in [-M, M]$ and  $\circ \in \{\neq, >, <\}$ , we have  $||f(a, b) - \mathbb{I}(a, b, \circ)||_{\infty} \leq \epsilon$ . 939

940 941 *Proof.* We first show that a Mamba block can implement  $\mathbb{I}[a > b]$  and  $\mathbb{I}[a < b]$ . For  $\mathbb{I}[a > b]$ , it is 942 equivalent to consider q(1, 0, a - b) where  $q(\cdot)$  defined in Lemma 4. So firstly we can use a linear layer with appropriate parameters  $W_0, b_0$  to convert the input [a, b] into the vector [1, 0, a-b]. Then 943 we can use Lemma 4 to implement  $\mathbb{I}[a > b]$  by changing the parameters of the first linear layer from 944  $\{W_1, b_1\}$  to  $\{W_1W_0, b_1 + W_1b_0\}$ . The proof for  $\mathbb{I}[a < b]$  is similar as well. 945

946 Noticing that  $\mathbb{I}[a \neq b] = 1 - (1 - \mathbb{I}[a > b]) \cdot (1 - \mathbb{I}[a < b])$ , we can implement  $\mathbb{I}[a \neq b]$ 947 through the following layers: Firstly, we can use one Mamba block to implement  $1 - \mathbb{I}[a > b]$  and 948  $1 - \mathbb{I}[a < b]$  simultaneously, where the hidden dimension will be 8 and the output is a vector  $[1 - \mathbb{I}[a < b]]$ 949  $\mathbb{I}[a > b], 1 - \mathbb{I}[a < b]]$ . Then, another Mamba block is constructed to implement the multiplication  $(1 - \mathbb{I}[a > b]) \cdot (1 - \mathbb{I}[a < b])$  according to Lemma 1 and the appropriate outermost linear layer 950 parameters are chosen to simultaneously achieve multiplication by a negative sign and addition of a 951 bias of 1, where the hidden dimension will be 4 and the output will be  $\mathbb{I}[a \neq b]$ . Thus, we complete 952 our proof. 953

954

926 927

931

Now, based on the basic operations that can be implemented by the Mamba blocks as discussed 955 above, we present the proof of Theorem 3: 956

957 **Theorem 7** (Perform DP problems with CoT). Considering any DP problem and given input sequences that satisfies Assumption 2, for any integer  $T \in \mathbb{N}$ , there exists several Mamba layers with 958 size O(T), such that the answer generated by the Mamba layers will be correct when the length of 959 the answer is no more than T. 960

961 962

963

*Proof.* Firstly, we illustrate the input format for the DP problem. We follow the embedding format in the proof of Theorem 4.7 in Feng et al. (2024), that is, assuming that the input at any step of solving the DP problem using CoT is a sequence of tokens embedded as follows:

$$\boldsymbol{x}_{t}^{(0)} = \left[\boldsymbol{e}_{t}^{\mathrm{input}}, \boldsymbol{e}_{t}^{\mathrm{state}}, \boldsymbol{e}_{t}^{\mathrm{dp}}, \boldsymbol{e}_{t}^{\mathrm{answer}}, \boldsymbol{e}_{t}^{\mathrm{sep}}, t, 1\right],$$

where the specific value of each part is depend on the content represented by the current token. More specifically, each part can be described as:

968 969

970 971

• If the current position denotes a input token, then we set  $e_t^{\text{input}}$  as the embedding of the input and simultaneously set  $e_t^{\text{state}} = e_t^{\text{dp}} = e_t^{\text{answer}} = e_t^{\text{sep}} = 0$ .

972 • If the current position is the final answer, then  $e_t^{\text{anwser}}$  denotes the embedding of the answer 973 and we set  $e_t^{\text{input}} = e_t^{\text{state}} = e_t^{\text{dp}} = e_t^{\text{sep}} = 0$ . 974 975 • If the current position denotes the j-th separator | between input sequences , then we set  $e_t^{\text{sep}} = e_t$  and  $e_t^{\text{input}} = e_t^{\text{state}} = e_t^{\text{dp}} = e_t^{\text{answer}} = 0$ . 976 977 • If the current position denotes an intermediate DP state, then we use  $e_t^{\text{state}}$  to denote the 978 embedding of the DP state and  $e_t^{dp}$  denotes the corresponding value. Similarly, other part 979 will be set to be **0**. 980 981 • The scalar t denotes the current position in the whole sequence, which holds the value for 982 all above cases. 983 984 We illustrate that here we use a concatenation operation to replace the residual connection defineed 985 in Eq (2), which is a technique also used by Feng et al. (2024); Yang et al. (2024) in similar proofs. 986 This is because, from the perspective of expressive capability, the two operations are equivalent: 987 the output of a Mamba block y = f(x) concatenated with the input (that is, |y, x|) can also be 988 represented using the residual connection: g([x, 0]) + [0, x] = [y, 0] + [0, x] = [y, x] where 989  $g: \mathbb{R}^{2d} \to \mathbb{R}^{2d}$  is another Mamba block and part of its parameters to be same as f and others 990 are set to be 0. Conversely, the concatenation can implement residual connection by using a linear 991 projection. 992 Here, we show our construction of several Mamba layers to solve the DP problem, which is com-993

posed of different blocks to perform different tasks:

- **Block 1:** The first block aims to calculate the problem size n and the embedding of the next state  $e_t^{\text{next\_state}}$ . This process can be described as follows:
  - Compute the problem size n: (i) First, we can replicate the position of the token  $t_{sep,1}, t_{sep,2}, \ldots, t_{sep,N}$  using the COPY operation. This can be achieved with a Mamba layer of size O(Ntd) according to Theorem 2; (ii) Then, we calculate the size of the problem as  $n = [t_{sep,1} 1, t_{sep,2} t_{sep,1} 1, \ldots, t_{sep,N} t_{sep,N-1} 1]$ , which can be done by applying a linear transformation using one Mamba layer, as shown by Lemma 3.
  - Obtain the next state  $e^{\text{next}$ -state: According to Assumption 2, the function  $e^{\text{next}$ -state} =  $f(n, e^{\text{state}})$  which determines the next state, can be approximated by constant-sized MLPs. Thus, this can also be implemented by having several Mamba layers degenerate into MLPs.
- The output after this step can be written as:

995

996 997 998

999

1000

1001

1002

1004

1012

1013

1014 1015

1016

1017

$$\boldsymbol{x}_t^{(1)} = [\boldsymbol{e}_t^{\text{input}}, \boldsymbol{e}_t^{\text{state}}, \boldsymbol{e}_t^{\text{next\_state}}, \boldsymbol{e}_t^{\text{dp}}, \boldsymbol{e}_t^{\text{answer}}, \boldsymbol{e}_t^{\text{sep}}, \boldsymbol{n}, t, 1]$$

**Block 2:** The second block is mainly constructed to find the indices of input tokens and intermediate DP values that are needed to calculate the DP value corresponding to  $e^{\text{next_state}}$ . Specifically, this can be described as follows:

- Calculate the needed indices: We calculate the positions of the input token  $p_t^s = I_s(n, e^{\text{state}})$  and the positions of tokens that correspond to needed DP values  $p_t^{\text{dp}} = I_{\text{dp}}(n, e^{\text{state}})$ . If  $I_s(n, e^{\text{state}}) = \emptyset$  or  $I_{\text{dp}}(n, e^{\text{state}}) = \emptyset$ , we set the positions to be a special value  $\gamma$ . According to Assumption 2, these two functions can be done by constant-size MLPs thus can be approximated by degenerated Mamba layers.
- Set the flag: (i) Set the flag  $f_t^{\text{answer}}$  based on whether the DP value of the current state is needed in the final aggregation function. This can be achieved by several Mamba layers with Assumption 2 that the function  $\mathcal{A} = f(n, s)$  can be approximated by MLPs and additionally using Lemma 5 to implement  $\mathcal{I}[e_t^{\text{state}} \neq e_j^{\text{state}}]$  where  $e_j^{\text{state}} \in \mathcal{A}$ . (ii) Set the flag  $f_t^{\text{state}}$  to denote whether the current state is the last state. This can be implemented by checking  $\mathbb{I}[e_t^{\text{next_state}} \neq 0]$  with Mamba layers using Lemma 5.

The output result after this step can be written as:

1028 1029

1032

1033

1034

1035

1039

1041

1043

1045

1046

1047

1050

1051

1056

1058

1062

1063

1064

$$\boldsymbol{x}_{t}^{(2)} = [\boldsymbol{e}_{t}^{ ext{input}}, \boldsymbol{e}_{t}^{ ext{state}}, \boldsymbol{e}_{t}^{ ext{nstate}}, \boldsymbol{e}_{t}^{ ext{dp}}, \boldsymbol{e}_{t}^{ ext{answer}}, \boldsymbol{e}_{t}^{ ext{sep}}, \boldsymbol{n}, \boldsymbol{p}_{t}^{s}, \boldsymbol{p}_{t}^{ ext{dp}}, f_{t}^{ ext{answer}}, f_{t}^{ ext{state}}, t, 1$$

Block 3: This block is designed to calculate the DP value for the next state. In detail, the implementation involves the following steps:

- Check the flag: We check the flag  $f_t^{\text{state}}$  using several Mamba layers using 5 to implement  $\mathbb{I}[f_t^{\text{state}} \neq 1]$  using Lemma 5. If  $f_t^{\text{state}} = 1$ , the current denote is the last state and we just need to set  $p_t = \gamma \mathbf{1}$  where  $p_t$  denotes  $p_t^s$  and  $p_t^{\text{dp}}$ , which implies  $I_s(n, e^{\text{state}}) = \emptyset$  and  $I_{\text{dp}}(n, e^{\text{state}}) = \emptyset$ , that is, no input tokens or DP values are needed.
  - Obtain the needed embeddings: If f<sub>t</sub><sup>state</sup> ≠ 1, then (i) We copy the input token embeddings e<sup>input</sup> at positions p<sub>t</sub><sup>s</sup>. This COPY operation can be implemented by a Mamba layer of size O(N<sub>s</sub>td) using Theorem 2; (ii) Simultaneously, we copy the embeddings of DP values at positions p<sub>t</sub><sup>dp</sup>, which can be achieved by a Mamba layer of size O(N<sub>dp</sub>td). If the position is empty, we just need to check I[p<sub>t</sub> ≠ γ1] and set the needed embeddings e<sup>input</sup> or e<sup>dp</sup> to be some special token. Totally, the size of Mamba layers in this step is O((N<sub>s</sub> + N<sub>dp</sub>)td).
  - Calculate the DP value: We calculate the DP value  $e_t^{\text{next_state}}$  for the next state with the Assumption 2 that the transition function can be approximated by several Mamba layers using Lemma 2.

# 1048 The output result after this step can be written as:

$$\boldsymbol{x}_t^{(2)} = [\boldsymbol{e}_t^{\text{input}}, \boldsymbol{e}^{\text{next\_state}}, \boldsymbol{e}_t^{\text{next\_dp}}, \boldsymbol{e}_t^{\text{answer}}, \boldsymbol{e}_t^{\text{sep}}, \boldsymbol{n}, f_t^{\text{answer}}, f_t^{\text{state}}, t, 1]$$

Block 4: The last block is constructed to implement the final aggregation function and output the final answer. Specifically, the steps are as follows:

- Check the flag: We identify whether the current state is the last state by checking  $\mathbb{I}[f_t^{\text{state}} \neq 1]$  by using Lemma 5. If  $f_t^{\text{state}} = 1$ , then all intermediate DP values have been solved and we need to compute the final answer.
- Obtain the needed embeddings: We collect the DP value embeddings  $e^{dp}$  of these tokens whose  $f^{answer} = 1$ , which can be achieved by COPY operation according to Theorem 2 with one Mamba layer of size  $O(N_A td)$ .
- Generate the final answer: Finally, we compute the answer by implementing the aggregation function, which can be achieved by constant-size MLPs according to Assumption 2, thus can also be achieved by several degenerated Mamba layers.

In summary, given a sequence length t and equipped with CoT, the parameter size required by the Mamba layers to generate the correct answer at each step is  $O(\tilde{N}td)$ , where  $\tilde{N} = \max\{N, N_s + N_{dp}, N_A\}$  is a constant independent of t, that is, the size of the Mamba layer scales linearly with t. Thus, we complete our proof.

1070 A.4 PROOF OF THEOREM 4

**Theorem 8** (Perform *m*-locality DP problems with CoT). Consider any *m*-locality DP problem and given input sequences that satisfies Assumption 2, for any integer  $T \in \mathbb{N}$ , there exists several Mamba layers with size O(m), such that the answer generated by the Mamba layers will be correct when the length of the answer is no more than T.

1076

**Proof.** The overall proof construction approach is similar to that of Theorem 3, with the only difference being that under the assumption of *m*-locality, when performing the COPY operation, the constructed Mamba only needs to focus on at most *m* tokens preceding the current position. This results in the size of the Mamba layers only needing to be  $O(\tilde{N}md)$ .

# 1080 A.5 More Details of Experiments

For the copy task experiments, we mainly refer to the setup by Jelassi et al. (2024). For Transformers, we select the GPT-NeoX architecture (Andonian et al., 2023) while for Mamba we use the Mamba GitHub repository(Gu & Dao, 2023) and all experiments are conducted on A800 GPUs.

More specifically, for the left part of Figure 2, we configure 10 layers for the Transformer (TF) and 5 layers for TF-small, with both having a hidden size of 1024 and RoPE (Su et al., 2024) as the positional encoding. For Mamba models, we configured 20 layers and a hidden size of 1024 for Mamba, 20 layers and a hidden size of 720 for Mamba-small-D, 10 layers and a hidden size of 1024 for Mamba-small-L, and 40 layers and a hidden size of 512 for Mamba-small-LD. We use an online sampling batch size of 8 and set the maximum context length to 220, meaning each example often contains multiple instances. AdamW(Loshchilov, 2017) is chosen as the optimizer with a learning rate of 1e-5 and weight decay of 0.1. We set  $N_{min} = 10$  and  $N_m ax = 30$  for all models. 

For the middle part of Figure 2, the Transformer and Mamba setups match the aforementioned configurations for TF and Mamba. Moreover, we set  $[N_{min}, N_{max}]$  to [5, 10], [10, 20], [20, 30] and [30, 40] for sequence lengths of 10, 20, 30 and 40 respectively.

For the right part of Figure 2, we use pretrained Mamba models of size 370M, 1.4B and 2.8B(Gu & Dao, 2023), which are pretrained on the Pile(Gao et al., 2020). For Mamba(2.8B)++, the prompt at the begining is just like:

"The following is a phonebook with the form: Gary Battle: 8444797678 Gary Gallegos: 9960330831.

1101 Remeber the phone number of Joseph Perry. Here is the phonebook:..."