

# STUFFED MAMBA: STATE EXPLOSION AND STATE CAPACITY OF RNN-BASED LONG-CONTEXT MODELING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

One essential advantage of recurrent neural networks (RNNs) over transformer-based language models is their linear computational complexity concerning the sequence length, which makes them much faster in handling long sequences during inference. However, most publicly available RNNs (e.g., Mamba and RWKV) are trained on sequences with less than 10K tokens, and their effectiveness in longer contexts remains largely unsatisfying so far. In this paper, we study the cause of the inability to process long context for RNNs and suggest critical mitigations. First, we investigate *state explosion* (SE) in Mamba-2 when processing long sequences, a phenomenon where some channels of the state exhibit exploding values that cause severe performance degradation. With controlled experiments, we discover that the model fails to forget the earlier tokens when there is more information than it can remember. We attribute this to overfitting due to the recurrent state being overparameterized for the training length, thereby establishing a relationship between SE and the capacity of the state. To support this hypothesis, we make an important empirical observation: for any given state size, there exists a training length threshold such that SE is exhibited if and only if the training length is greater than this threshold. Empirically searching for this threshold for different state sizes reveals that it is a linear function of the state size. We also search for the maximum context length at which the model can recall contextual information and find that this context length scales exponentially to the state size. Based on this, we empirically train a Mamba-2 370M with near-perfect passkey retrieval accuracy on 256K context length. This suggests a promising future for RNN-based long-context modeling. Code and model checkpoints will be publicly released.

## 1 INTRODUCTION

Recent transformer-based language models (Vaswani et al., 2017; Brown et al., 2020; Touvron et al., 2023; Achiam et al., 2023; Dubey et al., 2024) have demonstrated impressive capabilities in reasoning over long sequences with thousands and even millions of tokens (Team, 2024a;b; GLM et al., 2024). However, they rely on the attention mechanism that scales quadratically regarding the sequence length, making them extremely costly for inference over long sequences. In contrast, recurrent neural networks (RNNs) (Bengio et al., 1994) have a contextual memory with constant state size. Thus, during inference, their per-token computational and memory complexity scales linearly with the sequence length, making them much more efficient in processing long sequences.

Despite the great efficiency of RNNs in processing long contexts, their long-context performances are far from satisfying. Most recent state-of-the-art (SOTA) RNN-based language models (hereafter referred to as *RNNs* for simplicity), such as Mamba-1 (Gu and Dao, 2023), Mamba-2 (Dao and Gu, 2024), the RWKV series (Peng et al., 2023a; 2024), and GLA (Yang et al., 2023) are trained on sequences with less than 10K tokens. Existing works have shown that Mamba-1 and RWKV-4 suffer from severe performance drops when the context length exceeds their training length (Ben-Kish et al., 2024; Zhang et al., 2024a; Waleffe et al., 2024).

In this paper, we study the problem that what causes the current RNNs’ inability to handle long contexts and what are the possible solutions for supporting long contexts. When applying RNNs to longer contexts, we observe two critical problems. (1) RNNs are unable to generalize along the sequence length. They exhibit abnormal behavior when the context length exceeds the training length,

054 resulting in poor long-context performance. (2) Since their memory size is constant, although they  
055 can process infinitely long inputs, there is an upper bound to the amount of information the state  
056 can represent. Therefore, there is an upper bound of the contextual memory capacity (the maximum  
057 number of tokens that can be remembered), and tokens beyond that limit will be forgotten.

058 Then we dive deeper into the formation of the above problems. We first attribute the cause of  
059 the length generalization failure of SOTA RNNs to a phenomenon we call *state explosion* (SE). It  
060 describes the exploding channels in the recurrent state that cause performance degradation as the  
061 model consumes too many tokens. We inspect the memory state distribution over time and discover  
062 that the explosion is largely contributed by a few dominant outlier channels. These outliers cause  
063 vanishing values in other channels when the output hidden representation is normalized. By analyzing  
064 various components of the state update rule, we show that SE is caused by the inability to forget the  
065 earliest tokens (by decaying the state with a smaller multiplier) when there is more information than it  
066 can remember. To support this finding, we empirically show that it is possible to mitigate SE without  
067 training by forcing the model to forget contextual information by (1) reducing the memory retention  
068 and insertion strength and (2) reformulating the recurrence into an equivalent sliding window state  
069 (which makes it forget everything outside of the window). These methods allow Mamba-2 to consume  
070 more than 64K tokens without SE.

071 Additionally, we present the *state overparameterization hypothesis* as a possible explanation for why  
072 the model fails to learn a robust forgetting mechanism. This hypothesis tells us that one can avoid  
073 SE by training on sufficiently long sequences that exceed the model’s state capacity, which forces  
074 the model to learn how to forget excessive information to prevent SE. We empirically validate this  
075 hypothesis by discovering that for any state size, there is a threshold for training length beyond which  
076 the model will not exhibit SE. This insight allows us to establish a relationship between state capacity  
077 and state size. Then, by training Mamba-2 models of different sizes, we establish that the empirical  
078 capacity is a linear function of the state size. Furthermore, we conduct the same experiments on the  
079 widely used passkey retrieval task (Mohtashami and Jaggi, 2023), and show that the length where  
080 Mamba-2 has near-perfect passkey retrieval accuracy is an exponential function of the state size.  
081 The experiment results in a Mamba-2 370M model that can achieve near-perfect passkey retrieval  
082 accuracy on 256K context length, significantly outperforming transformer-based models of the same  
083 size in both retrieval accuracy and length generalizability. Our results show that the commonly  
084 used training lengths for RNN-based models may be suboptimal and that RNN-based long-context  
modeling has promising potential.

085 The main findings of this paper can be summarized as follows.

086 **State explosion:** In Mamba-2, there exists a small number of channels with exploding values that  
087 cause severe performance degradation on sequences longer than the training length. (Section 4)

088 **State overparameterization hypothesis:** By analyzing the hidden representations, we find that state  
089 explosion occurs because the model fails to forget when the state is being overflowed with excessive  
090 information. We attribute this to state overparameterization, which establishes the connection between  
091 SE and state capacity. (Section 4.3)

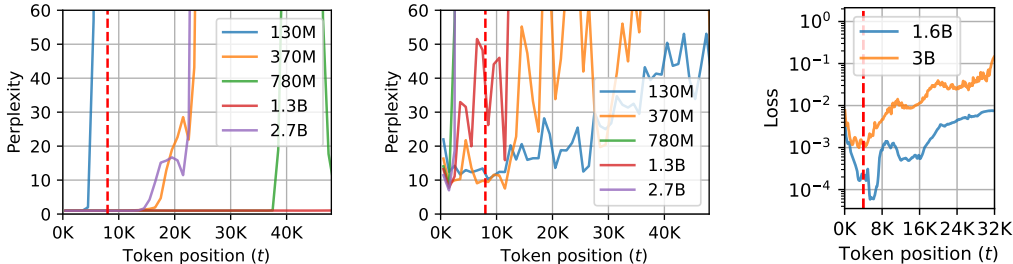
092 **Training length threshold:** In alignment with the state overparameterization hypothesis, we empiri-  
093 cally discover that for any state size, there exists a training length threshold where SE is exhibited if  
094 and only if the training length is below that threshold. (Section 5)

095 **Relationship between state capacity and state size:** Utilizing the previous finding, we can empiri-  
096 cally search for the relationship between the training length threshold and state size, and discover that  
097 it is linear. We also empirically search for state capacity on passkey retrieval, and find that it scaled  
098 exponentially with the state size. (Section 7.1)

## 101 2 RELATED WORKS

102 **RNN-Based Language Models** There is a recent surge of interest in RNN-based language models,  
103 because, contrary to transformer-based ones, their per-token inference cost does not increase with  
104 the sequence length. Linear Attention (Katharopoulos et al., 2020) replaces the softmax attention in  
105 transformer-based models with kernel-based approximations that have equivalent recurrent formula-  
106 tions. Some notable recent RNNs include the RWKV series (Peng et al., 2023a; 2024), the Mamba  
107

108  
109  
110  
111  
112  
113  
114  
115  
116  
117  
118  
119  
120  
121  
122  
123  
124  
125  
126  
127  
128  
129  
130  
131  
132  
133  
134  
135  
136  
137  
138  
139  
140  
141  
142  
143  
144  
145  
146  
147  
148  
149  
150  
151  
152  
153  
154  
155  
156  
157  
158  
159  
160  
161



(a) Mamba-2 Series on the “newlines” prompt (only “\n” characters). (b) Mamba-2 series on long documents from RedPajama. (c) RWKV-6 1.6B and 3B on the “newlines” prompt.

Figure 1: The perplexity of Mamba-2 and RWKV-6 as a function of token position on real and synthetic data. The red dotted line represents the models’ training lengths. They all fail to extrapolate.

series (Gu and Dao, 2023; Dao and Gu, 2024), Gated Linear Attention (Yang et al., 2023), among others (Zhang et al., 2024b; Yang et al., 2024; De et al., 2024; Arora et al., 2024; Orvieto et al., 2023; Sun et al., 2023). These models have shown strong capabilities in many language processing tasks, sometimes outperforming transformer-based models. However, as we will empirically show, some of these models fail to extrapolate much beyond their training length.

Some transformer-based models have adopted sliding window attention (Beltagy et al., 2020; Jiang et al., 2023), which essentially turns them into RNNs. However, these models have been shown to perform poorly in long-context tasks and fail to extrapolate to very long contexts (Zhang et al., 2024a).

**Length Generalization** Most SOTA language models in the last few years have been based on the transformer (Vaswani et al., 2017) architecture. These models, when using certain variants of position encoding, can process arbitrarily long sequences. However, they exhibit severe performance drops on tokens beyond the training length (Zhao et al., 2024). To alleviate this shortcoming, many works have focused on modifying positional encoding (Peng et al., 2023b; Zhu et al., 2023; Ding et al., 2024; Jin et al., 2024), some achieving training-free length generalization to certain extents. Similarly in spirit, this study also explores some post-training modifications for enhancing RNN-based models’ length generalization capabilities.

**Length Generalization of Mamba** Some prior works investigated the performance of Mamba as a function of context length (Park et al., 2024; Wen et al., 2024). Jelassi et al. (2024) empirically showed a sharp drop in performance beyond the training length for Mamba and LSTM on a copying task and also showed that Mamba struggles to copy from context unless its state size grows linearly with the context length. Arora et al. (2023) discussed the capacity and empirical ability to perform associative recall of transformer and some recurrent language models. In contrast, our paper explores the empirical relationship between state capacity and length generalization failure.

Some concurrent works have explored extending Mamba’s context length by controlling the discretization term ( $\Delta_t$  in Eq. 3) (Ben-Kish et al., 2024), such as dividing it by a constant to make it smaller<sup>1</sup>. This essentially makes the memory decay factor ( $\alpha_t$  in Eq. 5) closer to 1, which makes the state retain more contextual information. However, it also unnecessarily diminishes the inserted information on all tokens. Consequently, although it can mitigate SE, it results in poor performance in the passkey retrieval task (details in Appendix C).

### 3 PRELIMINARY

Most experiments in this study focus on Mamba-2 (Dao and Gu, 2024) because it has shown strong capabilities on several tasks and has publicly available checkpoints of multiple sizes, allowing us to

<sup>1</sup><https://github.com/jzhang38/LongMamba>

162 explore the relationship between state sizes and length limits. Moreover, it is more widely studied  
 163 than other RNNs, making it easier to use existing works as a reference.  
 164

165 **Mamba-2** The Mamba-2 architecture consists of  $L$  layers, each consisting of  $H$  heads computed in  
 166 parallel. The layer’s output is the sum of the heads’ outputs. Each head can be formulated as follows.  
 167

$$168 \quad y_t = \text{Norm}(o_t \odot u_t W_{\text{gate}}) W_o \in \mathbb{R}^d \quad (1)$$

$$169 \quad o_t = C_t h_t + D \odot x_t \in \mathbb{R}^P \quad (2)$$

$$170 \quad h_t = h_{t-1} \alpha_t + \bar{B}_t x_t \in \mathbb{R}^{N \times P} \quad (3)$$

172 where  $t$  denotes the current time step,  $u_t, y_t \in \mathbb{R}^d$  are the input and output hidden representations  
 173 of the  $t$ -th token,  $\text{Norm}(\cdot)$  denotes the RMS normalization (Zhang and Sennrich, 2019),  $D \in$   
 174  $\mathbb{R}^P$ ,  $W_{\text{gate}} \in \mathbb{R}^{d \times P}$ , and  $W_o \in \mathbb{R}^{P \times d}$  are trainable parameters,  $\odot$  denotes element-wise product,  
 175 and  $d, N, P$  are hyperparameters, denoting the hidden dimensionality, state dimension, and head  
 176 dimension, respectively. The other variables are parameterized as follows:

$$177 \quad \bar{B}_t = \Delta_t B_t \in \mathbb{R}^{N \times 1} \quad (4)$$

$$178 \quad \alpha_t = \exp(-\Delta_t \exp(A)) \in \mathbb{R} \quad (5)$$

$$179 \quad C_t = \sigma(\text{Conv}(u_t W_C)) \in \mathbb{R}^{1 \times N} \quad (6)$$

$$180 \quad B_t = \sigma(\text{Conv}(u_t W_B)) \in \mathbb{R}^{N \times 1} \quad (7)$$

$$181 \quad x_t = \sigma(\text{Conv}(u_t W_x)) \in \mathbb{R}^{1 \times P} \quad (8)$$

$$182 \quad \Delta_t = \text{Softplus}(u_t W_\Delta + b_\Delta) \in \mathbb{R} \quad (9)$$

185 where ( $W_C, W_B \in \mathbb{R}^{d \times N}$ ,  $W_x \in \mathbb{R}^{d \times P}$ ,  $W_\Delta \in \mathbb{R}^{d \times 1}$ ,  $b_\Delta, A \in \mathbb{R}$ ) are trainable model parameters.  
 186  $\text{Conv}(\cdot)$  denotes a channel-wise one-dimensional convolutional layer (see Appendix A.2).  $\sigma$  denotes  
 187 the SiLU function (Elfwing et al., 2017). Importantly,  $h_t$  is called the *hidden state* or *recurrent state*,  
 188 which is contextual memory that stores information from all tokens up to  $t$ .  $\alpha_t$  is the decay multiplier  
 189 that controls the strength of memory decay (i.e., forgetting). Appendix A presents more details.

190 It is worth mentioning that this update rule can be seen as a variant of Gated Linear Attention (Yang  
 191 et al., 2023) and is similar to many existing RNNs (e.g., RWKV (Peng et al., 2024) and RetNet (Sun  
 192 et al., 2023)). Thus, some conclusions/insights may apply to other architectures. We leave such  
 193 exhaustive ablation studies for future work.  
 194

## 195 4 STATE EXPLOSION

197 We first evaluate Mamba-2 and RWKV-6 on language modeling and passkey retrieval and find that  
 198 they fail to extrapolate beyond their training lengths. Then, we analyze the state’s distribution on  
 199 different context lengths and find that a few outlier channels exhibit exploding values. Based on this  
 200 finding, we coin the term *state explosion*, which describes the phenomenon where the recurrent state  
 201 exhibit exploding values, causing performance degradation.

202 By analyzing the components of the state’s update rule, we discover that SE is caused by the failure  
 203 to forget past information when the context is too long. We then propose a hypothesis that explains  
 204 SE from the perspective of overparameterization.  
 205

### 206 4.1 LENGTH GENERALIZATION FAILURE

208 **Language Modeling** Figure 1 shows the language modeling loss of Mamba-2 and RWKV-6 beyond  
 209 their training lengths. For controllability and to synthesize prompts of arbitrary lengths, this loss is  
 210 computed on a prompt consisting of only the “\n” characters, which we refer to as the “newlines”  
 211 prompt. However, we emphasize that the same observation also exists when processing texts from  
 212 the pre-training corpus. The result shows that both RNNs suffer great performance degradation when  
 213 the context length is much longer than their training lengths.  
 214

215 **Passkey Retrieval Evaluation** Language modeling may not reflect downstream capabilities, thus,  
 we evaluate several strong RNNs on the passkey retrieval task (Mohtashami and Jaggi, 2023; Zhang

216  
217  
218  
219  
220  
221  
222  
223  
224  
225  
226  
227  
228  
229  
230  
231  
232  
233  
234  
235  
236  
237  
238  
239  
240  
241  
242  
243  
244  
245  
246  
247  
248  
249  
250  
251  
252  
253  
254  
255  
256  
257  
258  
259  
260  
261  
262  
263  
264  
265  
266  
267  
268  
269

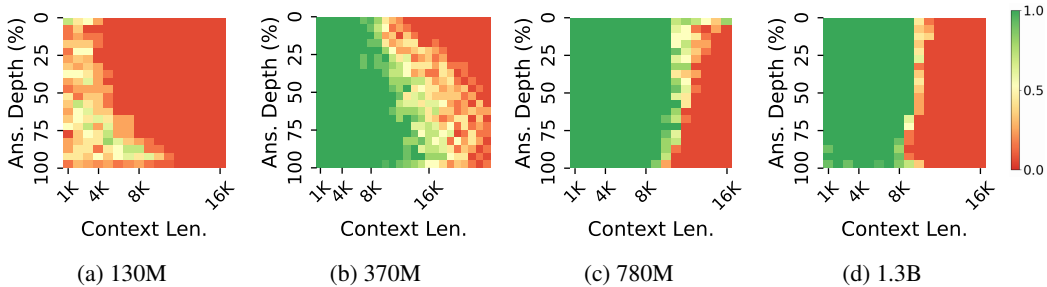


Figure 2: The performance of Mamba-2 official checkpoints on the passkey retrieval task. “Ans. Depth” refers to the position of the passkey relative to the context length, 0% means that it is located at the start of the context, 100% means that it is at the end.

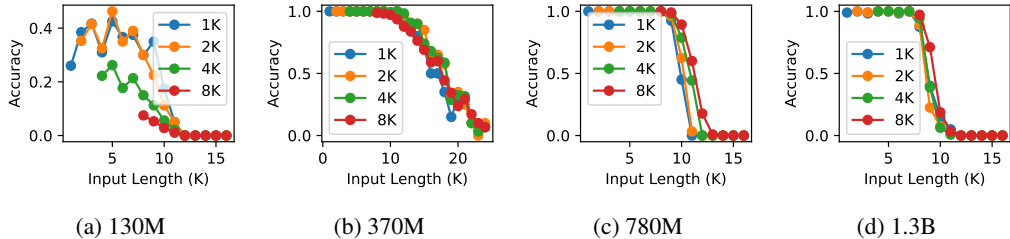


Figure 3: The accuracy of Mamba-2 official checkpoints on the passkey retrieval task where the answer is in the last  $n$  tokens, with  $n=1K, 2K, 4K,$  and  $8K$ .

et al., 2024a), a simple synthetic task where a model is prompted to recall a 5-digit *passkey* from a lengthy context. The hyperparameters and results of other RNNs can be found in Appendix B and D. The results for Mamba-2 are reported in Figure 2. We find that Mamba-2 models fail to generalize to sequences longer than the training length. For instance, Mamba-2, trained on 8K context windows, has near-perfect retrieval accuracy within 8K contexts (except for the smaller 130M checkpoint), but poor or even zero accuracy on sequences longer than 16K, regardless of model sizes.

This behavior is unexpected because the update rule (Eq. 3) has a stable exponential memory decay (it converges to a constant value if the variables are fixed). Therefore, we expect RNNs of such form to have a good retrieval accuracy on the last  $k$  tokens, and tokens earlier than that are forgotten. This unexpected finding also implies that when processing contexts longer than the training length  $T_{train}$ , it is better to keep just the last  $T_{train}$  tokens and discard everything else. However, this is not trivial in an online inference scenario because all token information is compressed into a single state.

#### 4.2 OBSERVATION OF STATE EXPLOSION

Since the recurrent state’s dimensionality does not change over time, the sharp change of behavior during length generalization must be a result of a change in the state’s value. We inspect the statistics of the recurrent states of each layer in Mamba-2 370M<sup>2</sup> and find that the mean and variance of some heads change sharply when the context length exceeds the training length. One example is shown in Figure 4<sup>3</sup>. Appendix G reports the statistics of each layer. The state at  $t = 20K$  of one head with exploding variance is shown in Figure 5. From it, we discover that this variance explosion can be largely attributed to a few outlier channels, while most channels are relatively stable.

We emphasize that **SE occurs largely independent of the prompt**, occurring in both pre-training data samples and generated meaningless texts, even for prompts consisting of only whitespace characters

<sup>2</sup>We are using this model size as an example, but the same observation is found with any model size.  
<sup>3</sup>We emphasize that we are showing only a subset of heads for one layer in this figure to keep it clean. Exploding heads are found in other layers as well.

270  
271  
272  
273  
274  
275  
276  
277  
278  
279

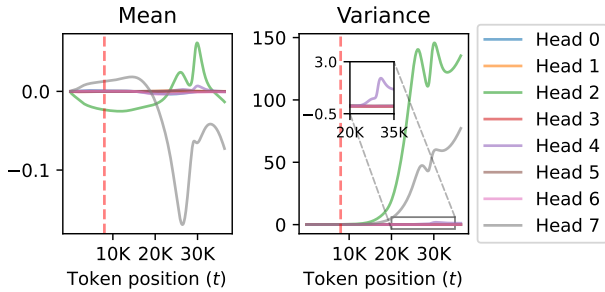


Figure 4: The mean and variance of the first 8 heads in the 38th layer of Mamba-2 370M. It exhibits a clear explosion when  $t$  is greater than the training length. The red dotted line indicates the training length.

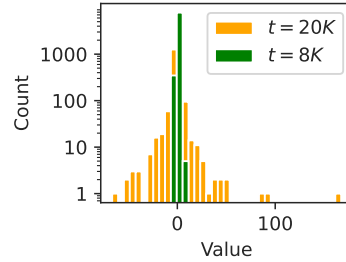
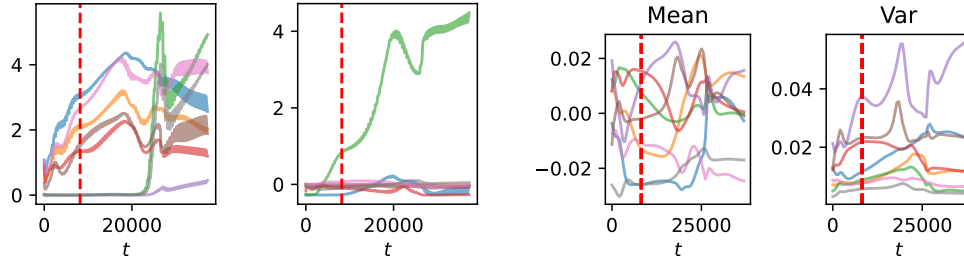


Figure 5: The distribution of the channels in a collapsing state (head 2 for the 38th layer) at two different time steps.

280  
281  
282  
283  
284  
285  
286  
287  
288  
289  
290  
291  
292  
293  
294



(a)  $\Delta_t$ . Each curve represents a head. (b)  $B_t$ . Each curve represents a channel. (c) The statistics of  $x_t$  over time. Each curve represents a head.

Figure 6: The value of the components in the update rule ( $\Delta_t$ ,  $B_t$ , and  $x_t$ ) on some heads with SE in the 38th layer in Mamba-2 370M. The red dotted line indicates the training length.

295  
296  
297

(Figure 1 (a) and (b) shows SE on two different prompts). This means that the information inserted into the state does not cause its explosion.

300  
301  
302  
303  
304  
305  
306  
307  
308  
309

To further attribute to the specific variables that cause SE, we inspect the values of  $\Delta_t$ ,  $B_t$ , and  $x_t$  on various heads with exploding states. Figure 6 reports one example of the inspection, we can see that  $x_t$  is relatively stable compared to the  $\Delta_t$  and  $B_t$ , even though they are all functions of  $u_t$ . We also notice that  $B_t$  explodes earlier than  $\Delta_t$  (ignoring the non-exploding heads). Further inspection reveals that the convolutional weights that generate  $\Delta_t$  and  $B_t$  (Eq. 7 and 9) are noticeably greater in variance than those for  $x_t$  (Eq. 8). A more in-depth attribution study is outside the scope of this work.

### 310 4.3 STATE OVERPARAMETERIZATION HYPOTHESIS

311  
312  
313  
314  
315

Here, we present a high-level explanation for SE. We argue that **SE arises from state overparameterization** relative to the training length. In other words, the state is excessively large for the training length, allowing the model to achieve strong language modeling performance without learning how to forget when the state is about to overflow.

316  
317

To support this argument, we formulate the hidden state as a weighted sum of previously inserted information:

318  
319  
320  
321

$$h_t = \sum_{i=1}^t \alpha_{i:t} \bar{B}_i x_i, \quad \alpha_{i:t} = \left( \prod_{j=i}^t \alpha_j \right) \in (0, 1) \quad (10)$$

322  
323

Thus,  $\alpha_{i:t}$  describes the memory strength about the  $i$ -th token at  $t$  time step. Figure 7 shows the memory strength of the first token at different time steps, and we find that the exploded heads (heads 2, 4, and 7 in the 38th layer) have a strong inclination toward retaining all information within the

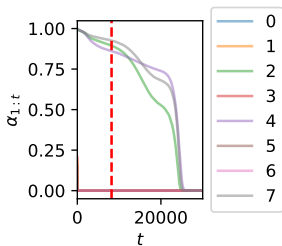


Figure 7: The cumulative decay of the first token ( $\alpha_{1,t}$ ) over time. Each line represents a head. The statistics are collected from the 38th layer of Mamba-2 370M.

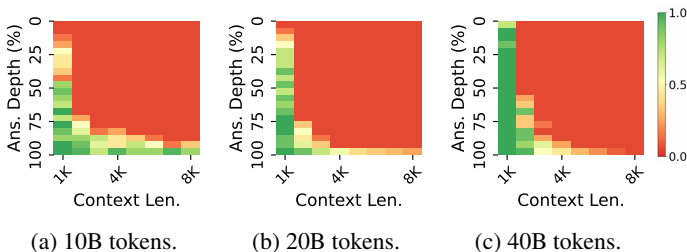


Figure 8: The passkey retrieval results of intermediate checkpoints during the pre-training of Mamba-2 370M on 512 sequence length. SE only occurs in the model beyond a certain amount of training data.

training length, with a memory strength of over 0.8 at  $t=8K$ . This implies that the model has not learned to forget information (by producing a smaller decay  $\alpha_j$ ), but it still has decent language modeling capabilities because the state has the capacity to store all information in 8K tokens.

Furthermore, we pre-train Mamba-2 from scratch with  $T_{\text{train}} = 512$  and evaluate the intermediate checkpoints on passkey retrieval, as reported in Figure 8. It shows that SE is only exhibited by checkpoints beyond a certain amount of training, which coincides with behaviors of overfitting—a result of overparameterization. One can also notice that the overfitted checkpoint outperforms earlier checkpoints on shorter sequences, which further strengthens the hypothesis that the model converges to less forgetting.

Finally, in Section 7.1, we empirically find that for any given state size, there exists a length threshold  $T_*$  where SE is exhibited if and only if  $T_{\text{train}} < T_*$ .

#### 4.4 HOW TO MITIGATE STATE EXPLOSION?

To further support the state overparameterization hypothesis in the previous section, we demonstrate that it is possible to mitigate SE by simply forcing the model to forget more and/or remember less information. This can be achieved by training on sufficiently long sequences or by training-free modifications to the update rule. We emphasize that the training-free methods serve as demonstrations of the link between forgetting and SE, and we recommend training on longer sequences in practice.

##### 4.4.1 TRAINING ON LONGER SEQUENCES

We employ two techniques to improve effectiveness and efficiency when training on longer sequences.

**Data Engineering** To ensure that the data contains as much long-term structure as possible, we filter out sequences with less than 4K tokens. Buckman and Gelada (2024) have shown that this is critical for training effective long-context models. Although we train on sequences longer than 4K tokens, we do not use a higher length threshold because the above threshold already removes about 97.6% of the data in the original corpus. To train on longer sequences, we simply concatenate sequences and delimit them with a special EOS (End-of-Sequence) token.

**Truncated Backpropagation Through Time** In the vanilla Mamba-2, the states are initialized to zeros for each data sample. Instead, we initialize the states as the final state of the previous sequence. This is equivalent to concatenating multiple sequences, but stopping the backpropagation of gradients at certain intervals. This technique has been shown to help extend the context length of RNNs (Yang et al., 2023) and alleviate the memory cost of caching activations for computing gradients. Based on Yang et al. (2023) and our preliminary tests, we use concatenate 12 sequences with this technique by default.

#### 4.4.2 TRAINING-FREE MITIGATION METHODS

Here, we propose some tricks for mitigating SE by inducing more forgetting without training. Since the performance of these tricks is not the focus of this paper, we report their evaluation result in Appendix I. The main takeaway from these methods is that SE can be prevented by simply inducing more forgetting by modifying the update rule without additional training.

**Reduced Retention and Insertion Strength (RRI)** This method assumes that  $\alpha_t$  and  $B_t$  control the memory retention and insertion strength, respectively. We scale them with a multiplier smaller than 1. The actual value is chosen by validation on the “newlines” prompt. Existing works (?) have experimented with modifying  $\Delta_t$ , but it controls both the insertion and decay strength, making it hard to analyze and control.

**Sliding Window** We can utilize the fact that the state  $h_t$  can be written as a weighted sum (Eq. 10) to simulate a sliding window mechanism without re-processing from the start of the window at every step. Let  $w \in \mathbb{N}$  denote the window size and  $h_t^{(r)} \in \mathbb{R}^{N \times P}$  denote the hidden state when applying the model on the last  $w$  tokens at time step  $t$ . We can then compute  $h_t^{(r)}$  exactly as the difference between two states:

$$h_t^{(r)} = \sum_{i=t-r+1}^t \alpha_{i:t} \bar{B}_i x_i = \sum_{i=1}^t \alpha_{i:t} \bar{B}_i x_i - \alpha_{t-r+1:t} \sum_{i=1}^{t-r} \alpha_{i:t-r} \bar{B}_i x_i = h_t - \alpha_{t-r+1:t} h_{t-r} \quad (11)$$

During streaming generation, we only have to maintain  $(h_{t-1}, h_{t-r}, \alpha_{t-r+1:t})^4$ , and advance each of them in parallel. However, directly computing  $\alpha_{t:t-r}$  may suffer from instability due to floating-point imprecision. Therefore, we maintain  $\Delta_{t-r:t} = \sum_{i=t-r}^t \Delta_i$  instead, and re-compute  $\alpha_{t-r:t} = \exp(-\Delta_{t-r:t} \exp(A))$  at every step, which incurs minimal computational cost.

This method can be applied to all RNNs that can be written as a weighted sum, which includes RWKV 5 and 6, RetNet, GLA, etc. It doubles the computation and memory cost for generation, but we believe that it is an acceptable trade-off because RNNs have a very low generation cost compared to transformer-based models and the context processing cost is unchanged.

## 5 STATE CAPACITY

Based on the discussion and hypothesis in Section 4.3, we can avoid SE by training on sufficiently long sequences such that the amount of contextual information exceeds the capacity of the state, thereby forcing the model to learn how to forget. This hypothesis implies the following law:

Let  $P_S$  and  $T_{\text{train}}$  denote the recurrent state size and training length, respectively, there exists a training length threshold  $T_*(P_S)$  such that SE is exhibited if and only if  $T_{\text{train}} < T_*$ .

In Section 7.1, we empirically validated this by sweeping different training lengths for different state sizes, and checking whether the model exhibits SE. Hence,  $T_*$  can be viewed as an indirect measurement of state capacity because it tells us when the amount of contextual information exceeds the capacity. Here, it is important to disambiguate two notions of state capacity: the theoretical capacity of the state and the empirically observed capacity on a given task. In this paper, we refer to the latter, which is an empirical description of the amount of context information a state can store without forgetting.

We conduct the training as described in Section 4.4.1. We train multiple Mamba-2 with different state sizes and training lengths to find the relationship between  $T_*$  and  $P_S$ . To determine whether a state has exploded, we feed the “newlines” prompt with 1M tokens to the model, and define explosion as the point where perplexity is more than 2x the maximum perplexity within  $T_{\text{train}}$  tokens.

<sup>4</sup>We also have to cache the last  $r$  token IDs, but their size is negligible compared to  $h_{t-1}$  and  $h_{t-r}$ .



432  
433  
434  
435  
436  
437  
438  
439  
440  
441  
442  
443  
444  
445  
446  
447  
448  
449  
450  
451  
452  
453  
454  
455  
456  
457  
458  
459  
460  
461  
462  
463  
464  
465  
466  
467  
468  
469  
470  
471  
472  
473  
474  
475  
476  
477  
478  
479  
480  
481  
482  
483  
484  
485

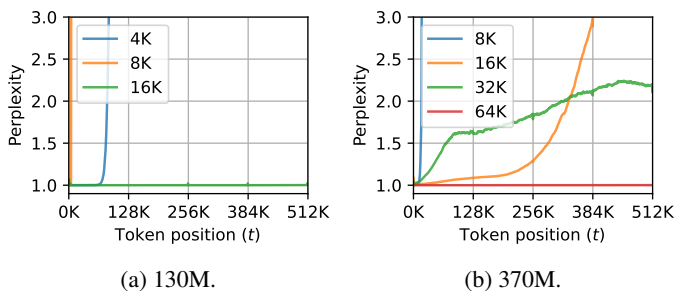


Figure 9: The perplexity as a function of token position for Mamba-2 on the “newlines” prompt with different training lengths.

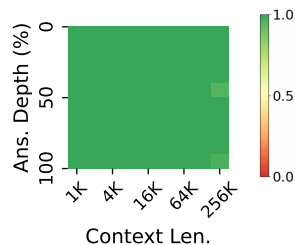


Figure 10: Passkey retrieval performance of Mamba-2 370M. The x-axis is in the log scale.

## 5.1 STATE CAPACITY IN PASSKEY RETRIEVAL

It is worth highlighting that the model does not necessarily fail to recall information beyond the last  $T_*$  tokens, especially when there is a clear distinction between the target information and other contextual information. Therefore, we also search for the state capacity on the passkey retrieval task, which we regard as an evaluation of the minimum ability to accurately recall a short sequence from a lengthy context (the task is very easy compared to other long-context tasks). Similar to the previous section, we train with different lengths for different state sizes and identify the maximum context length where the model has an accuracy over 95%, which we regard as the *state capacity in passkey retrieval*, denoted with  $C_{\text{passkey}}$ . In this task, the noisy context is repetitive, thus, the amount of contextual information is largely independent of the context length, therefore, the capacity should grow roughly exponentially with the state size.<sup>5</sup>

## 6 EXPERIMENTS

We briefly describe the experimental details for training with longer sequences. See Appendix F for more comprehensive experimental details.

**Data** We start from RedPajama-V2 (Computer, 2023), an open dataset with 30T tokens extracted from CommonCrawl<sup>6</sup>, and we perform deduplication to ensure data quality. During evaluation, we sample documents longer than 16K tokens and concatenate them if it is not long enough.

**Models** We experiment with six model configurations with different state sizes to find the relationship between state capacity and size. For each of them, we perform an extensive search with training lengths up to 256K tokens. To save cost, we continue pre-train from three official checkpoints of Mamba-2 of size 130M, 370M, and 780M. They were pre-trained with 8K sequences. The other three model configurations (36M, 47M, and 85M) are trained from scratch.

**Hyperparameters** We use the WSD (Hu et al., 2024) with 10% decay steps. This scheduler is chosen because it is competitive with the commonly used cosine scheduler while allowing simple resumption from intermediate checkpoints, saving large amounts of computational resources. We report the result of the best checkpoint selection by validation on passkey retrieval.

<sup>5</sup>If we train Mamba-2 on passkey retrieval data, the model can theoretically handle infinitely long contexts. Here, the model is only trained with the next token prediction objective, which means the model will *not* ignore the irrelevant context, and the ability to retain information for extended time emerges from language modeling.

<sup>6</sup><https://commoncrawl.org/>

486  
487  
488  
489  
490  
491  
492  
493  
494  
495  
496  
497  
498  
499  
500  
501  
502  
503  
504  
505  
506  
507  
508  
509  
510  
511  
512  
513  
514  
515  
516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527  
528  
529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539

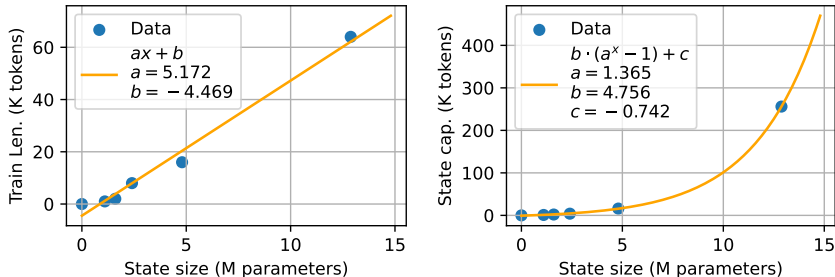


Figure 11: The left figure shows the training length beyond which the model does not exhibit SE (i.e.,  $T_*$  in Section 5). The right figure shows the state capacity on passkey retrieval (i.e.,  $C_{\text{passkey}}$ ) as a function of state size.

## 7 RESULTS

### 7.1 LENGTH GENERALIZATION BY TRAINING ON LONGER SEQUENCES

In Figure 9, we plot the language modeling perplexity as a function of token position for Mamba-2 130M and 370M with different training lengths. We can see that for each model size, there is a training length threshold, beyond which the model has much better length extrapolation, which supports our arguments discussed in Section 4.4.1.

### 7.2 STATE CAPACITY AS A FUNCTION OF STATE SIZE

Figure 11 shows the state capacity of Mamba-2 on language modeling and passkey retrieval. The rightmost data point in both plots corresponds to Mamba-2 370M. We have confirmed that the 780M model (with a state size of 19.3M) also exhibits SE at training lengths below 128K, but do not have enough resources to train the model beyond this length. The results establish a linear relationship  $T_* = 5.172 \cdot P_S - 4.469$  between the length  $T_{\text{train}} = T_*$  at which SE stops occurring and the state size  $P_S$ . The  $R^2$  value is over 0.999. This indicates that **to train a Mamba-2 with robust length generalization, one should use training lengths that grow linearly with the state size.**

The second plot of Figure 11 shows that the capacity of Mamba-2 on passkey retrieval is exponential concerning the state size, the function is  $C_{\text{passkey}} = 4.756 \cdot (1.365^{P_S} - 1) - 0.742$ , with an  $R^2$  value over 0.999. This is because the amount of information in the context does not increase with its length. In other words, we are storing a constant amount of information while the number of combinations of the state grows exponentially with the number of elements. Figure 10 shows the best checkpoint of Mamba-2 370M on passkey retrieval. The result is very promising because, to the best of our knowledge, no previous models with less than 1B model parameters have near-perfect accuracy at this length in this task.

## 8 CONCLUSION

This paper discovers and presents the first systemic study on *state explosion* (SE), a phenomenon in RNNs that causes length generalization failure. Controlled experiments on Mamba-2 reveals that SE occurs because the model fails to forget when the state is about to overflow, and we conclude that this phenomenon is caused by an overparameterized state with excessive state capacity. Then we show that SE can be mitigated by inducing more forgetting by modifying the update rule or by training on context lengths that exceed the *state capacity*. We discover that for any state size, there exists a training length beyond which SE will not occur. With this insight, we empirically estimate the state capacity of Mamba-2 on language modeling and the passkey retrieval task. With some simple data engineering and state initialization tricks, we achieve much better performance with Mamba-2 on the passkey retrieval task than existing models. Our results indicate that Mamba-2 not only is highly efficient in handling long sequences but also has great performance potential.

## LIMITATIONS

All models studied in this work can be seen as a specific case of *linear attention* models, whose recurrent state is decayed by an element-wise or scalar gate (they can be viewed as variants of Gated Linear Attention (Yang et al., 2023)). We have chosen to study these models because of their strong capabilities, yet, some conclusions may not be directly transferred to other variants of RNNs.

Our continued pre-training approach for extending the context length of RNNs can be rather expensive, some of the models require training with up to 50B tokens, which is 1/6 of their pre-training amount of data. Also, to ensure simplicity, controllability, and generality, we have not used more advanced techniques for training long-context models, such as upsampling longer data samples, better data order or format, using data with more long-distance dependencies, etc.

The passkey retrieval task that we used extensively in this study is very simple. Hence, high accuracy on this task may not reflect the capabilities of the model on real-world long-context tasks, because that requires more advanced capabilities such as high-resolution retrieval, reasoning, state-tracking, etc. The result is nonetheless promising because it indicates that the model can recall the correct information and further capabilities may be achieved by building on the recalled information.

SE is somewhat prompt-dependent. While we found that for models with greatly overparameterized states, SE is highly consistent across different prompts, certain models with less overparameterization may exhibit SE on some prompts while successfully extrapolating indefinitely on others. Attributing SE to specific features of the prompt is a promising future research direction. Moreover, as evident in Figure 1, length generalization failure is less severe in RWKV-6. This is because RWKV-6 has a much smaller state, as shown in Figure 12. Some architectural choices may also affect the occurrence of SE, but we leave such an exploration for future work.

## REFERENCES

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Simran Arora, Sabri Eyuboglu, Aman Timalsina, Isys Johnson, Michael Poli, James Zou, Atri Rudra, and Christopher Ré. 2023. [Zoology: Measuring and improving recall in efficient language models](#). *Preprint*, arXiv:2312.04927.
- Simran Arora, Sabri Eyuboglu, Michael Zhang, Aman Timalsina, Silas Alberti, Dylan Zinsley, James Zou, Atri Rudra, and Christopher Re. 2024. Simple linear attention language models balance the recall-throughput tradeoff. In *Workshop on Efficient Systems for Foundation Models II @ ICML2024*, volume abs/2402.18668.
- Iz Beltagy, Matthew E Peters, and Arman Cohan. 2020. Longformer: The long-document transformer. *arXiv preprint arXiv:2004.05150*.
- Assaf Ben-Kish, Itamar Zimmerman, Shady Abu-Hussein, Nadav Cohen, Amir Globerson, Lior Wolf, and Raja Giryes. 2024. [Decimamba: Exploring the length extrapolation potential of mamba](#). *Preprint*, arXiv:2406.14528.
- Yoshua Bengio, Patrice Simard, and Palo Frasconi. 1994. Learning long-term dependencies with gradient descent is difficult. *IEEE transactions on neural networks*, 5(2):157–166.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language models are few-shot learners](#). *Preprint*, arXiv:2005.14165.
- Jacob Buckman and Carles Gelada. 2024. Compute-optimal Context Size.
- Together Computer. 2023. [Redpajama: an open dataset for training large language models](#).

- 594 Tri Dao and Albert Gu. 2024. Transformers are ssms: Generalized models and efficient algorithms  
595 through structured state space duality. *arXiv preprint arXiv:2405.21060*.  
596
- 597 Soham De, Samuel L. Smith, Anushan Fernando, Aleksandar Botev, George-Cristian Muraru, Albert  
598 Gu, Ruba Haroun, Leonard Berrada, Yutian Chen, Srivatsan Srinivasan, Guillaume Desjardins,  
599 Arnaud Doucet, David Budden, Yee Whye Teh, Razvan Pascanu, Nando de Freitas, and Caglar  
600 Gulcehre. 2024. Griffin: Mixing Gated Linear Recurrences with Local Attention for Efficient  
601 Language Models. *arXiv*, abs/2402.19427.
- 602 Yiran Ding, Li Lina Zhang, Chengruidong Zhang, Yuanyuan Xu, Ning Shang, Jiahang Xu, Fan Yang,  
603 and Mao Yang. 2024. Longrope: Extending llm context window beyond 2 million tokens. *arXiv  
604 preprint arXiv:2402.13753*.
- 605 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha  
606 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of  
607 models. *arXiv preprint arXiv:2407.21783*.  
608
- 609 Stefan Elfving, Eiji Uchibe, and Kenji Doya. 2017. [Sigmoid-weighted linear units for neural network  
610 function approximation in reinforcement learning](#). *Preprint*, arXiv:1702.03118.
- 611 Team GLM, Aohan Zeng, Bin Xu, Bowen Wang, Chenhui Zhang, Da Yin, Diego Rojas, Guanyu  
612 Feng, Hanlin Zhao, Hanyu Lai, Hao Yu, Hongning Wang, Jiada Sun, Jiajie Zhang, Jiale Cheng,  
613 Jiayi Gui, Jie Tang, Jing Zhang, Juanzi Li, Lei Zhao, Lindong Wu, Lucien Zhong, Mingdao Liu,  
614 Minlie Huang, Peng Zhang, Qinkai Zheng, Rui Lu, Shuaiqi Duan, Shudan Zhang, Shulin Cao,  
615 Shuxun Yang, Weng Lam Tam, Wenyi Zhao, Xiao Liu, Xiao Xia, Xiaohan Zhang, Xiaotao Gu,  
616 Xin Lv, Xinghan Liu, Xinyi Liu, Xinyue Yang, Xixuan Song, Xunkai Zhang, Yifan An, Yifan Xu,  
617 Yilin Niu, Yuantao Yang, Yueyan Li, Yushi Bai, Yuxiao Dong, Zehan Qi, Zhaoyu Wang, Zhen  
618 Yang, Zhengxiao Du, Zhenyu Hou, and Zihan Wang. 2024. [Chatglm: A family of large language  
619 models from glm-130b to glm-4 all tools](#). *Preprint*, arXiv:2406.12793.
- 620 Albert Gu and Tri Dao. 2023. [Mamba: Linear-time sequence modeling with selective state spaces](#).  
621 *ArXiv*, abs/2312.00752.  
622
- 623 Shengding Hu, Yuge Tu, Xu Han, Chaoqun He, Ganqu Cui, Xiang Long, Zhi Zheng, Yewei Fang,  
624 Yuxiang Huang, Weilin Zhao, et al. 2024. Minicpm: Unveiling the potential of small language  
625 models with scalable training strategies. *arXiv preprint arXiv:2404.06395*.
- 626 Samy Jelassi, David Brandfonbrener, Sham M. Kakade, and Eran Malach. 2024. [Repeat after me:  
627 Transformers are better than state space models at copying](#). *Preprint*, arXiv:2402.01032.  
628
- 629 Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot,  
630 Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al.  
631 2023. Mistral 7b. *arXiv preprint arXiv:2310.06825*.
- 632 Hongye Jin, Xiaotian Han, Jingfeng Yang, Zhimeng Jiang, Zirui Liu, Chia-Yuan Chang, Huiyuan  
633 Chen, and Xia Hu. 2024. Llm maybe longlm: Self-extend llm context window without tuning.  
634 *arXiv preprint arXiv:2401.01325*.
- 635 Angelos Katharopoulos, Apoorv Vyas, Nikolaos Pappas, and François Fleuret. 2020. Transformers  
636 are rnns: Fast autoregressive transformers with linear attention. In *International conference on  
637 machine learning*, pages 5156–5165. PMLR.
- 638 Amirkeivan Mohtashami and Martin Jaggi. 2023. Landmark attention: Random-access infinite  
639 context length for transformers. In *Workshop on Efficient Systems for Foundation Models@  
640 ICML2023*.  
641
- 642 Antonio Orvieto, Samuel L. Smith, Albert Gu, Anushan Fernando, Caglar Gulcehre, Razvan Pascanu,  
643 and Soham De. 2023. [Resurrecting recurrent neural networks for long sequences](#). *Preprint*,  
644 arXiv:2303.06349.  
645
- 646 Jongho Park, Jaeseung Park, Zheyang Xiong, Nayoung Lee, Jaewoong Cho, Samet Oymak, Kang-  
647 wook Lee, and Dimitris Papailiopoulos. 2024. [Can mamba learn how to learn? a comparative  
study on in-context learning tasks](#). *Preprint*, arXiv:2402.04248.

- 648 Bo Peng, Eric Alcaide, Quentin G. Anthony, Alon Albalak, Samuel Arcadinho, Stella Biderman,  
649 Huanqi Cao, Xin Cheng, Michael Chung, Matteo Grella, G Kranthikiran, Xuming He, Haowen  
650 Hou, Przemyslaw Kazienko, Jan Kocoń, Jiaming Kong, Bartłomiej Koptyra, Hayden Lau, Krishna  
651 Sri Ipsit Mantri, Ferdinand Mom, Atsushi Saito, Xiangru Tang, Bolun Wang, Johan Sokrates Wind,  
652 Stansilaw Wozniak, Ruichong Zhang, Zhenyuan Zhang, Qihang Zhao, Peng Zhou, Jian Zhu, and  
653 Rui Zhu. 2023a. [Rwkv: Reinventing rnns for the transformer era](#). In *Conference on Empirical  
654 Methods in Natural Language Processing*.
- 655 Bo Peng, Daniel Goldstein, Quentin Anthony, Alon Albalak, Eric Alcaide, Stella Biderman, Eugene  
656 Cheah, Xingjian Du, Teddy Ferdinan, Haowen Hou, Przemysław Kazienko, Kranthi Kiran GV,  
657 Jan Kocoń, Bartłomiej Koptyra, Satyapriya Krishna, Ronald McClelland Jr., Niklas Muennighoff,  
658 Fares Obeid, Atsushi Saito, Guangyu Song, Haoqin Tu, Stanisław Woźniak, Ruichong Zhang,  
659 Bingchen Zhao, Qihang Zhao, Peng Zhou, Jian Zhu, and Rui-Jie Zhu. 2024. Eagle and Finch:  
660 Rwkv with Matrix-Valued States and Dynamic Recurrence. *arXiv*, abs/2404.05892.
- 661  
662 Bowen Peng, Jeffrey Quesnelle, Honglu Fan, and Enrico Shippole. 2023b. Yarn: Efficient context  
663 window extension of large language models. *arXiv preprint arXiv:2309.00071*.
- 664 Zhen Qin, Songlin Yang, Weixuan Sun, Xuyang Shen, Dong Li, Weigao Sun, and Yiran Zhong. 2024.  
665 Hgrn2: Gated Linear RNNs with State Expansion. In *First Conference on Language Modeling*,  
666 volume abs/2404.07904.
- 667  
668 Yutao Sun, Li Dong, Shaohan Huang, Shuming Ma, Yuqing Xia, Jilong Xue, Jianyong Wang, and  
669 Furu Wei. 2023. Retentive Network: A Successor to Transformer for Large Language Models.  
670 *arXiv*, abs/2307.08621.
- 671  
672 Gemini Team. 2024a. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of  
673 context. *arXiv*.
- 674  
675 Magic Team. 2024b. [100m token context windows](#).
- 676  
677 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée  
678 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand  
679 Joulin, Edouard Grave, and Guillaume Lample. 2023. [Llama: Open and efficient foundation  
680 language models](#). *Preprint*, arXiv:2302.13971.
- 681  
682 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz  
683 Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information  
684 processing systems*, 30.
- 685  
686 Roger Waleffe, Wonmin Byeon, Duncan Riach, Brandon Norick, Vijay Korthikanti, Tri Dao, Albert  
687 Gu, Ali Hatamizadeh, Sudhakar Singh, Deepak Narayanan, Garvit Kulshreshtha, Vartika Singh,  
688 Jared Casper, Jan Kautz, Mohammad Shoeybi, and Bryan Catanzaro. 2024. [An empirical study of  
689 mamba-based language models](#). *Preprint*, arXiv:2406.07887.
- 690  
691 Kaiyue Wen, Xingyu Dang, and Kaifeng Lyu. 2024. [Rnns are not transformers \(yet\): The key  
692 bottleneck on in-context retrieval](#). *Preprint*, arXiv:2402.18510.
- 693  
694 Songlin Yang, Bailin Wang, Yikang Shen, Rameswar Panda, and Yoon Kim. 2023. Gated linear  
695 attention transformers with hardware-efficient training. *arXiv preprint arXiv:2312.06635*.
- 696  
697 Songlin Yang, Bailin Wang, Yu Zhang, Yikang Shen, and Yoon Kim. 2024. Parallelizing Linear  
698 Transformers with the Delta Rule over Sequence Length. *arXiv*, abs/2406.06484.
- 699  
700 Biao Zhang and Rico Sennrich. 2019. [Root mean square layer normalization](#). *Preprint*,  
701 arXiv:1910.07467.
- 702  
703 Xinrong Zhang, Yingfa Chen, Shengding Hu, Zihang Xu, Junhao Chen, Moo Hao, Xu Han, Zhen  
704 Thai, Shuo Wang, Zhiyuan Liu, and Maosong Sun. 2024a. [∞Bench: Extending long context  
705 evaluation beyond 100K tokens](#). In *Proceedings of the 62nd Annual Meeting of the Association  
706 for Computational Linguistics (Volume 1: Long Papers)*, pages 15262–15277, Bangkok, Thailand.  
707 Association for Computational Linguistics.

702 Yu Zhang, Songlin Yang, Ruijie Zhu, Yue Zhang, Leyang Cui, Yiqiao Wang, Bolun Wang, Freda  
703 Shi, Bailin Wang, Wei Bi, Peng Zhou, and Guohong Fu. 2024b. [Gated slot attention for efficient](#)  
704 [linear-time sequence modeling](#). *Preprint*, arXiv:2409.07146.

705  
706 Liang Zhao, Xiaocheng Feng, Xiachong Feng, Dongliang Xu, Qing Yang, Hongtao Liu, Bing Qin,  
707 and Ting Liu. 2024. [Length extrapolation of transformers: A survey from the perspective of](#)  
708 [positional encoding](#). *Preprint*, arXiv:2312.17044.

709 Dawei Zhu, Nan Yang, Liang Wang, Yifan Song, Wenhao Wu, Furu Wei, and Sujian Li. 2023.  
710 Pose: Efficient context window extension of llms via positional skip-wise training. *arXiv preprint*  
711 *arXiv:2309.10400*.

712  
713  
714  
715  
716  
717  
718  
719  
720  
721  
722  
723  
724  
725  
726  
727  
728  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755

## A MAMBA-2 ARCHITECTURE

For completeness, we give a more detailed formulation of the Mamba-2 architecture here, although we recommend the readers refer to the original paper (Dao and Gu, 2024) or a detailed blog post by the authors<sup>7</sup>. The model accepts a sequence of  $T$  token IDs as input  $\mathbf{I} = [i_1, \dots, i_T] \in \mathbb{R}^T$ ,  $i_t \in \{1, 2, \dots, V\}$ , where  $V$  denotes the vocabulary size. It performs next token prediction by predicting the probability distribution over the vocabulary as each time step, denoted as  $P \in \mathbb{R}^{T \times V}$ . The model can be formulated as follows.

$$\begin{aligned} \mathbf{U}^{(0)} &= \text{Embed}_{\text{in}}(\mathbf{I}) \in \mathbb{R}^{T \times d} \\ \mathbf{U}^{(l)} &= \text{Mamba}^{(l)} \left( \text{Norm} \left[ \mathbf{U}^{(l-1)} \right] \right) \in \mathbb{R}^{T \times d} \\ \mathbf{P} &= \text{Embed}_{\text{out}} \left( \text{Norm} \left[ \mathbf{U}^{(L)} \right] \right) \in \mathbb{R}^{T \times V} \end{aligned}$$

where  $L$  denotes the number of layers,  $l \in \{1, \dots, L\}$  denotes the layer index,  $\mathbf{U}^{(l)} \in \mathbb{R}^{T \times d}$  represents the input of the  $l$ -th layer,  $\mathbf{U}^{(0)}$  represents the input of the first layer.  $\text{Mamba}^{(l)}(\cdot)$  denotes the  $l$ -th Mamba layer,  $\text{Embed}_{\text{in}}(\cdot)$  and  $\text{Embed}_{\text{out}}(\cdot)$  denote the input and output embedding layers, and  $\text{Norm}(\cdot)$  denotes RMS normalization (Zhang and Sennrich, 2019).  $d$  denotes the number of dimensions of each token embedding. Similar to many other models, Mamba-2 ties the weight of the input and output embedding layers.

Each Mamba layer consists of  $H$  ‘‘heads’’ that are computed in parallel. The result of which is summed together. The  $t$ -th token ( $t \in \{1, \dots, T\}$ ) in a head is computed as follows.

$$\begin{aligned} y_t &= \text{Norm}(o_t \odot u_t W_{\text{gate}}) W_o \in \mathbb{R}^d \\ o_t &= C_t h_t + D \odot x_t \in \mathbb{R}^P \\ h_t &= h_{t-1} \exp(-\Delta_t \exp(A)) + \Delta_t B_t^T x_t \in \mathbb{R}^{N \times P} \\ C_t &= \sigma(\text{Conv}(u_t W_C)) \in \mathbb{R}^N \\ B_t &= \sigma(\text{Conv}(u_t W_B)) \in \mathbb{R}^N \\ x_t &= \sigma(\text{Conv}(u_t W_x)) \in \mathbb{R}^P \\ \Delta_t &= \text{Softplus}(u_t W_\Delta + b_\Delta) \in \mathbb{R} \end{aligned}$$

$u_t$  denotes the  $t$ -th input representation. In other words, for the  $l$ -th layer, we have  $\mathbf{U}^{(l)} = [u_1^{(l)}, \dots, u_T^{(l)}]$ ,  $\mathbf{U}^{(l+1)} = [y_1^{(l)}, \dots, y_T^{(l)}]$ , and  $u_t^{(l+1)} = y_t^{(l)}$ .  $\text{Conv}(\cdot)$  denotes a channel-wise one-dimensional convolutional layer with a kernel size of 4, and  $\sigma$  denotes the SiLU activation function (Elfwing et al., 2017). All vectors are row vectors, and the result of a matrix multiplied by a scalar is the matrix with each element multiplied by that scalar.

( $W_{\text{gate}}, W_x \in \mathbb{R}^{d \times P}$ ,  $W_o \in \mathbb{R}^{P \times d}$ ,  $W_C, W_B \in \mathbb{R}^{d \times N}$ ,  $W_\Delta \in \mathbb{R}^{d \times 1}$ ,  $b_\Delta, A \in \mathbb{R}$ ) are trainable parameters of the layer, and  $P, N$  are hyperparameters. The authors call  $P$  the *head dimension* and  $N$  the *state size*. In practice, the weights of  $W_B, W_C$  are shared among different heads.

### A.1 STATE SIZE

The authors of Mamba-2 always set  $P = 64$ ,  $N = 128$ , and  $H = 2d/P$ . Thus, the state size of each Mamba-2 layer is  $HPN = 2dN = 256d$ . In transformer-based models, when using multi-headed attention, usually, the product of the head count  $H$  and head dimension  $P$  equals the hidden dimension  $d$ . Therefore, the KV cache of a transformer-based model is  $2Td$ , which means that when using the same hidden dimension, the state of a Mamba-2 layer is equal in size to a KV cache of 128 tokens.

Compared to many other recurrent models (e.g., the RWKV series (Peng et al., 2023a; 2024), GLA (Yang et al., 2023), and RetNet (Sun et al., 2023)), Mamba-2 does not have a state-less feed-forward network and has considerably more heads in each layer, making the state size much larger than other recurrent models. Compared to Mamba-1 (Gu and Dao, 2023), Mamba-1 uses  $N = 16$ , which means that the state size in Mamba-2 is 8 times larger than the state in a Mamba-1 model of roughly the model parameter count. Figure 12 shows the relationship between state size and model size of the RNN models in this study.

<sup>7</sup><https://tridao.me/blog/2024/mamba2-part1-model/>

810  
811  
812  
813  
814  
815  
816  
817  
818  
819  
820  
821  
822  
823  
824  
825  
826  
827  
828  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
860  
861  
862  
863

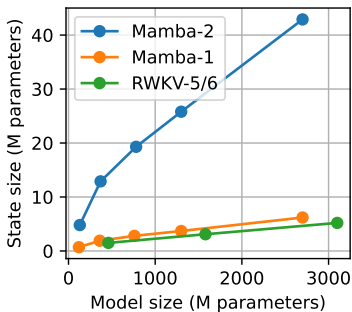


Figure 12: The relationship between state size and model size of various RNN models considered in this paper.

### A.2 SHORT CONVOLUTION

The  $\text{Conv}(\cdot)$  function in Mamba-2 is a one-dimensional convolutional layer applied to each channel separately. For  $i$ -th channel, it can be formulated as follows.

$$y_{t,i} = \sum_{j=1}^k w_{j,i} x_{t-j,i} \in \mathbb{R}, \quad i = 1, \dots, n_c$$

$k$  denotes the kernel size (set to 4 by default).  $i$  denotes the channel index,  $n_c$  denotes the number of channels.  $y_{t,i} \in \mathbb{R}$  denotes the  $i$ -th channel of the output vector at  $t$ -th time step.  $x_{t,i}$  represents the  $i$ -th channel of the input vector at  $t$ -th time step.  $w_{j,i} \in \mathbb{R}$  denotes the  $j$ -th value in the convolutional kernel for channel  $i$ .

This model component accepts the last 4 token embeddings at the input. Therefore, it also has a state that contains information about the context, which we refer to as the *convolutional state*. To be concrete, due to information propagation through the layers, the short convolutional layer is a function of the last  $4L$  tokens. For the 370M model size, this length is  $4 \times 48 = 192$ . Therefore, we can reasonably assume that this component contains much less contextual information relative to the recurrent state  $h_t$ . Thus, we have largely ignored this state in various discussions in this paper. However, we have also reported the distribution of the input to this short convolutional layer over time in Figure 17, for reference. As we can see, the convolutional state is relatively stable over time (compared to the recurrent state).

## B PASKEY RETRIEVAL INFERENCE PARAMETERS

Throughout the whole paper, we use greedy decoding, not just for reproducibility, but also because our preliminary results show that other decoding parameters give noticeably worse performance on passkey retrieval.

We use 32-bit floating point precision for both model parameters and activations during inference, to ensure that precision errors do not introduce noise to the result. We have conducted some preliminary evaluations with BF16 and FP16 and found that there are no noticeable differences with using FP16, but computing some activations, especially the  $\Delta_t$  and  $\alpha_t$  with BF16 introduces an error around  $1e-3$ . However, the explosion of channels in the states is consistently observed despite this precision error.

### B.1 PASKEY RETRIEVAL PROMPT

The prompt that we use for the passkey retrieval task is as follows, using 34847 as the passkey for example, which is adapted from existing works (Zhang et al., 2024a). We also evaluate with slight variations to the template in preliminary experiments but do not observe considerable differences in the results.

There is important info hidden inside a lot of irrelevant text.  
Find it and memorize it.



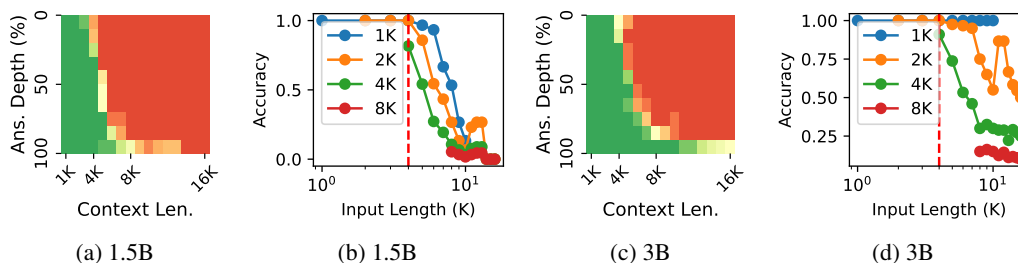
864 The grass is green. The sky is blue. The sun is yellow. Here we  
 865 go. There and back again.  
 866 ...  
 867 The grass is green. The sky is blue. The sun is yellow. Here we  
 868 go. There and back again.  
 869 The passkey is 34847. Remember it. 34847 is the passkey.  
 870 The grass is green. The sky is blue. The sun is yellow. Here we  
 871 go. There and back again.  
 872 ...  
 873 The grass is green. The sky is blue. The sun is yellow. Here we  
 874 go. There and back again.  
 875 What is the passkey? The passkey is

877 We sweep different context lengths  $T \in \{1K, 2K, \dots, 256K\}$ , and for each length  $T$ , we generate  
 878  $n$  prompts with evenly distributed needle positions, i.e., the  $i$ -th needle ( $i \in \{0, \dots, n-1\}$ ) of a  
 879 sample is inserted at position  $T \times i/n - 1$ , from the beginning.

## 882 C MAMBA-2 WITH MODIFIED $\Delta_t$ ON PASSKEY RETRIEVAL

884 Ben-Kish et al. (2024) and GitHub user jzhang28<sup>8</sup> propose to improve Mamba’s length generalization  
 885 by reducing the value of  $\Delta_t$ . Ben-Kish et al. (2024) propose a heuristic method for identifying which  
 886 head to modify and how to modify  $\Delta_t$ . However, their method requires task-dependent tweaking, so  
 887 we do not consider comparing against it. jzhang28 propose to simply multiply  $\Delta_t$  by a constant (they  
 888 used 0.5). We apply this method and sweep different  $\Delta_t$  for the best passkey retrieval performance,  
 889 but they all result in worse performance than the original model across all context lengths.

## 891 D PASKEY RETRIEVAL EVALUATION WITH OTHER ARCHITECTURES



904 Figure 13: The performance of RWKV-5 official checkpoints on the passkey retrieval task. Each  
 905 curve in (b) and (d) represents the accuracy of retrieving the needle when it is within the last  $r$  tokens,  
 906 with  $r \in \{1K, 2K, 4K, 8K\}$ .

908 Here, we also evaluate RWKV-5, RWKV-6, and Mamba-1 (some popular and strong RNNs) on the  
 909 passkey retrieval task. The result is reported in Figure 13, 14 and 15. We can see that SE is observed  
 910 in Mamba-1, but it is less severe for RWKV-5 and RWKV-6. We hypothesize that this difference is a  
 911 result of architectural differences and that the state size is smaller in RWKV-5 and RWKV-6.

## 913 E PRE-TRAINED CHECKPOINTS

914 The pre-trained checkpoints used in our experiments are given in Table 1.

917 <sup>8</sup><https://www.github.com/jzhang38/LongMamba>

918  
919  
920  
921  
922  
923  
924  
925

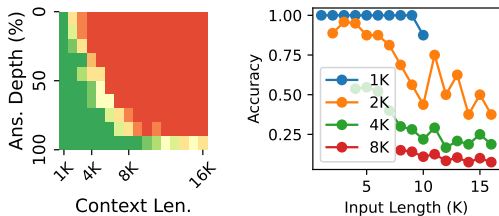


Figure 14: RWKV-6 1.6B result on the passkey retrieval task. The left plot shows the retrieval accuracy of the needle when it appears in the last  $r = \{1K, 2K, 4K, 8K\}$  tokens.

930  
931  
932  
933  
934  
935  
936  
937  
938  
939  
940

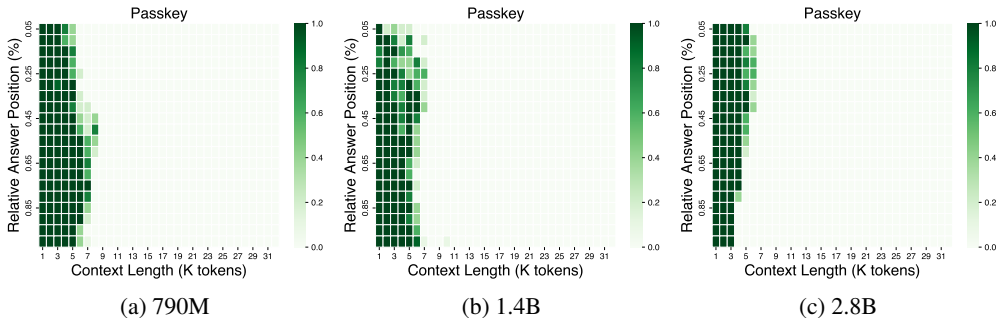


Figure 15: The performance of Mamba-1 official checkpoints on the Passkey task. We can see a clear exhibition of state explosion, similar to Mamba-2.

944  
945  
946  
947  
948  
949  
950  
951  
952  
953  
954  
955  
956  
957

Model	Checkpoint URLs
RWKV-5	<a href="https://huggingface.co/RWKV/rwkv-5-world-all-ptn">https://huggingface.co/RWKV/rwkv-5-world-all-ptn</a>
RWKV-6	<a href="https://huggingface.co/RWKV/v6-Finch-1B6-HF">https://huggingface.co/RWKV/v6-Finch-1B6-HF</a> <a href="https://huggingface.co/RWKV/v6-Finch-3B-HF">https://huggingface.co/RWKV/v6-Finch-3B-HF</a>
Mamba-1	<a href="https://huggingface.co/state-spaces/mamba-130m">https://huggingface.co/state-spaces/mamba-130m</a> <a href="https://huggingface.co/state-spaces/mamba-370m">https://huggingface.co/state-spaces/mamba-370m</a> <a href="https://huggingface.co/state-spaces/mamba-790m">https://huggingface.co/state-spaces/mamba-790m</a> <a href="https://huggingface.co/state-spaces/mamba-1.4b">https://huggingface.co/state-spaces/mamba-1.4b</a> <a href="https://huggingface.co/state-spaces/mamba-2.8b">https://huggingface.co/state-spaces/mamba-2.8b</a>
Mamba-2	<a href="https://huggingface.co/state-spaces/mamba2-130m">https://huggingface.co/state-spaces/mamba2-130m</a> <a href="https://huggingface.co/state-spaces/mamba2-370m">https://huggingface.co/state-spaces/mamba2-370m</a> <a href="https://huggingface.co/state-spaces/mamba2-780m">https://huggingface.co/state-spaces/mamba2-780m</a> <a href="https://huggingface.co/state-spaces/mamba2-1.3b">https://huggingface.co/state-spaces/mamba2-1.3b</a> <a href="https://huggingface.co/state-spaces/mamba2-2.7b">https://huggingface.co/state-spaces/mamba2-2.7b</a>

Table 1: The pre-trained checkpoints used in our experiments.

961

F EXPERIMENTAL DETAILS OF TRAINING ON LONGER SEQUENCES

962  
963  
964  
965  
966  
967  
968  
969  
970  
971

We perform a hyperparameter search on learning rates, sweeping  $\{1e-5, 2e-5, 5e-5, 1e-4, 2e-4, 5e-4, 1e-3\}$ , selecting the best performing one by validation on passkey retrieval<sup>9</sup>. Regarding the WSD scheduler, it warms up linearly for 1000 steps and decays linearly with 50K steps. This setup is inspired by the authors of WSD (Hu et al., 2024).

Other hyperparameters are kept as similar to the original papers for Mamba-2 as possible. That means we use 0.5M tokens per batch because we found this to give more stable results for continual pre-training instead of the 1M batch size from the original paper. Training is done mainly in BF16,

<sup>9</sup>While the loss of many checkpoints was highly similar, their performance in passkey retrieval can differ a lot.

with some activations in FP32 (in the same manner as the official implementation). The optimizer is AdamW, with a 0.1 weight decay. Moreover, we use 1.0 gradient clipping.

All experiments are run on A800 80G, some are run with multiple nodes, and others with multiple GPUs on a single node.

## F.1 MODEL CONFIGURATIONS

For the models smaller than the 130M official checkpoint, we pre-train from scratch using the configurations reported in Table 2. We try to follow the same depth-to-width ratio found in the official checkpoints, although the ratio is not entirely consistent in those checkpoints. Hyperparameters not mentioned are kept the same as the 130M checkpoint.

Model size	State size	# Layers	Hidden dim.	# heads
<i>Official checkpoints</i>				
780M	19.3M	48	1536	48
370M	12.9M	48	1024	32
130M	4.8M	24	768	24
<i>Our checkpoints trained from scratch</i>				
84.6M	2.4M	12	768	24
47.0M	1.6M	12	512	16
36.4M	0.8M	6	512	16

Table 2: The configurations of the models used in finding the passkey retrieval memory capacity as a function of the state size.

## G STATE STATISTICS OVER CONTEXT LENGTH

Here, we provide a more detailed result on the inspection of SE over time.

Figure 16 shows the hidden state of the recurrent mechanism described in Eq. 3. Additionally,  $B_t$ ,  $C_t$ , and  $x_t$  in Mamba-2 are generated with a short channel-wise convolutional layer with a kernel size of 4:

$$\begin{aligned} B_t &= \sigma(\text{Conv}[u_t W_B]) \\ C_t &= \sigma(\text{Conv}[u_t W_C]) \\ x_t &= \sigma(\text{Conv}[u_t W_x]) \end{aligned}$$

where  $\sigma$  is the SiLU activation function. This function is also stateful because it operates on the last 4 tokens, therefore, we also collect the statistics of this convolutional state and report them in Figure 17. As we can see, the convolutional states are much more stable compared to the recurrent states. This is because only the last 4 tokens contribute to this state which avoids the explosion as a result of cumulative sum.

## H HGRN-2 LENGTH GENERALIZATION

We additionally evaluate HGRN-2 (Qin et al., 2024) on the newlines prompt, and find that it also exhibits severe performance degradation on the “newlines” prompt. The model size is 1.3B. Perhaps surprisingly, the increase in perplexity happens considerably before the context length reaches the training length.

## I TRAINING-FREE LENGTH GENERALIZATION RESULT

Figure 19 reports the result of the training-free length generalization methods on Mamba-2 780M, evaluated on concatenated documents from RedPajama. We can see that while LongMamba<sup>10</sup> can

<sup>10</sup><https://github.com/jzhang38/LongMamba>

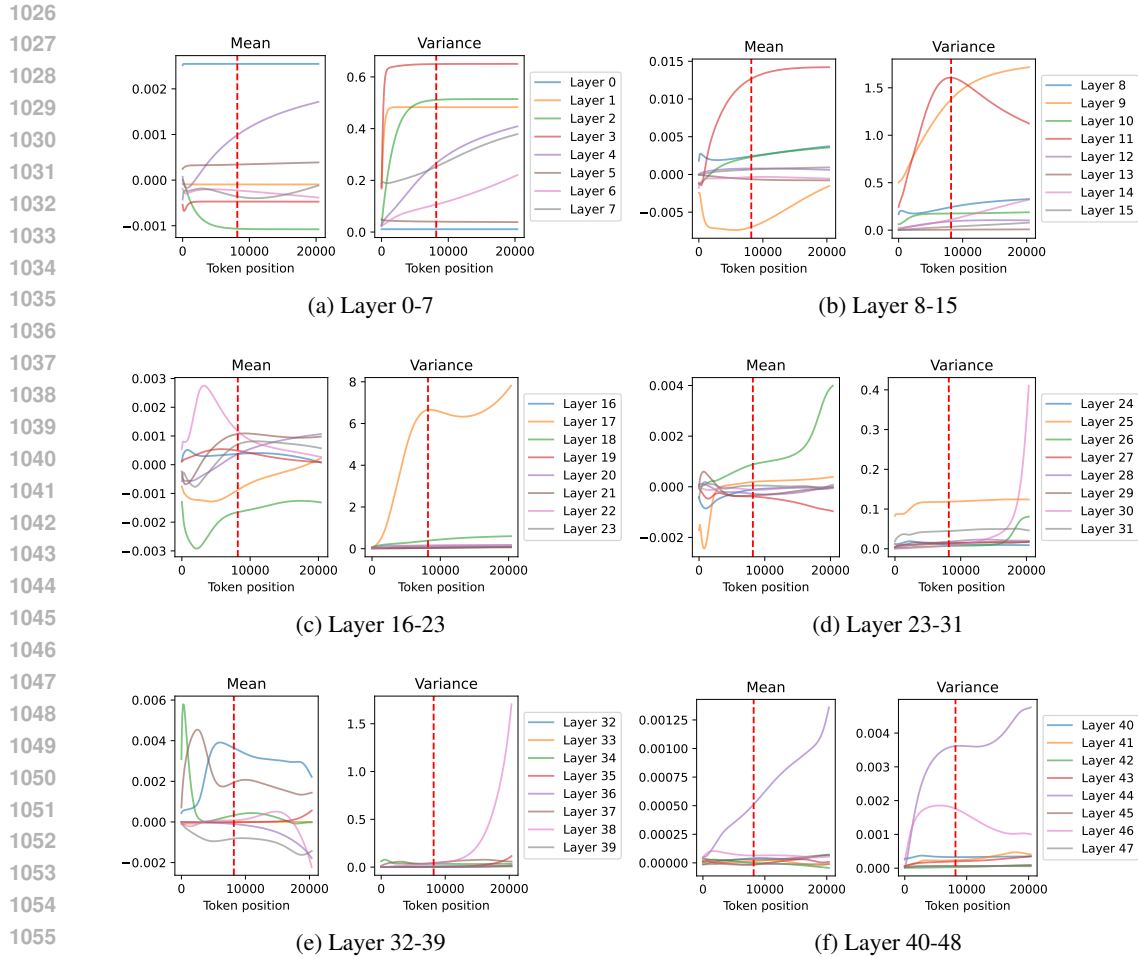


Figure 16: The mean and variance of the hidden state of each layer of Mamba-2 370M, computed on the “newlines” prompt.

greatly improve the length generalizability of the model by more than 3x, it causes noticeably greater perplexity on shorter sequences and still inevitably exhibits SE. All our methods successfully suppress SE, allowing the model to generalize to more than 64K tokens, although they underperform LongMamba on a range of context length (roughly from 8K to 20K).

## J THE “NEWLINES” PROMPT

In this paper, we collect the statistics of the state computed on a “newlines” prompt, a prompt where every token is the newline token, as shown below.

```
\n\n\n\n\n\n\n\n\n\n\n\n\n\n\n\n...
```

However, we again emphasize that SE is observed on prompts extracted from the pre-training corpus, the passkey retrieval task, or other randomly generated sequences. We have chosen the “newlines” prompt because the samples from the pre-training corpus are too short, and this prompt produces the most consistent and smooth layer statistics.

1080  
1081  
1082  
1083  
1084  
1085  
1086  
1087  
1088  
1089  
1090  
1091  
1092  
1093  
1094  
1095  
1096  
1097  
1098  
1099  
1100  
1101  
1102  
1103  
1104  
1105  
1106  
1107  
1108  
1109  
1110  
1111  
1112  
1113  
1114  
1115  
1116  
1117  
1118  
1119  
1120  
1121  
1122  
1123  
1124  
1125  
1126  
1127  
1128  
1129  
1130  
1131  
1132  
1133

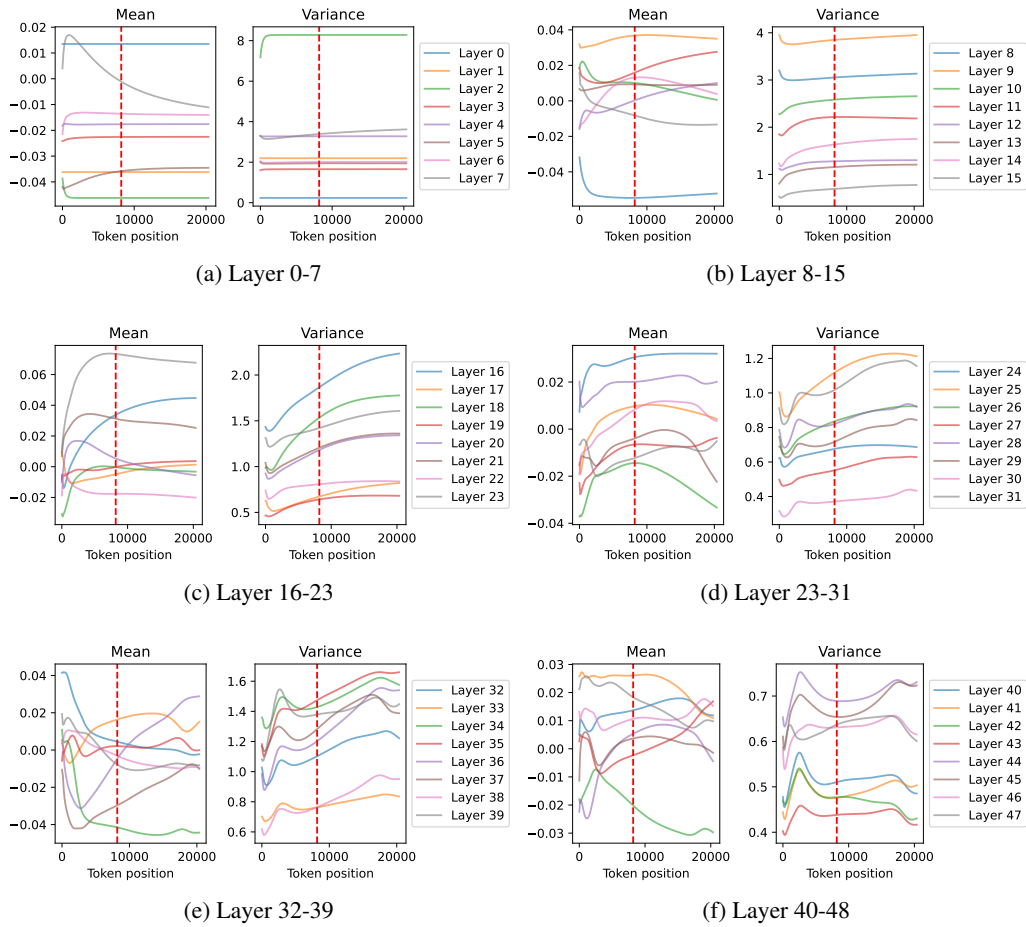


Figure 17: The mean and variance of the convolutional states (the representation of the last four tokens) of each layer in Mamba-2 370M, computed on the “newlines” prompt. We can see that the mean and variance are visibly more stable than the recurrent state.

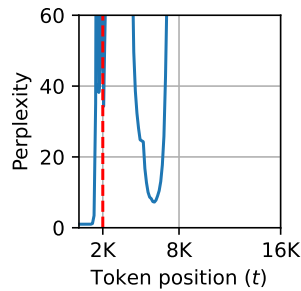


Figure 18: The perplexity of HGRN-2 1.3B on the “newlines” prompt as a function of time.

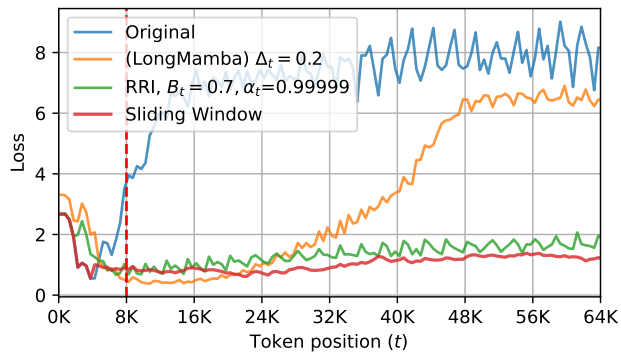


Figure 19: Result of training-free length generalization methods described in Section 4.4.2. The loss is computed on long documents from RedPajama (Computer, 2023). The red dotted line represents the training length.