

---

# Mamba State-Space Models Can Be Strong Downstream Learners

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Mamba [22] state-space models (SSMs) have recently outperformed state-of-the-  
2 art (SOTA) Transformer large language models (LLMs) in various tasks and been  
3 widely adapted. However, Mamba’s downstream learning capabilities remain ei-  
4 ther unexplored—e.g., mixed-precision (MPFT) and parameter-efficient fine-tuning  
5 (PEFT)—or under-evaluated—e.g., in-context learning (ICL). For the latter, recent  
6 works [45, 19] reported Mamba’s ICL rivals SOTA Transformer LLMs using non-  
7 standard benchmarks. In contrast, we show that on standard benchmarks, pretrained  
8 Mamba models achieve only 38% of the ICL performance improvements (over  
9 zero-shot) of comparable Transformers.

10 Enabling MPFT and PEFT in Mamba architectures is challenging due to recurrent  
11 dynamics and highly customized CUDA kernels, respectively. However, we prove  
12 that Mamba’s recurrent dynamics are robust to small input changes using dynamical  
13 systems theory. Empirically, we show that performance changes in Mamba’s  
14 inference and fine-tuning due to mixed-precision align with Transformer LLMs.  
15 Furthermore, we show that targeting key memory buffers in Mamba’s customized  
16 CUDA kernels for low-rank adaptation regularizes SSM parameters, thus achieving  
17 parameter efficiency while retaining speedups. We show that combining MPFT and  
18 PEFT enables up to 2.15 times more tokens-per-second and 65.5% reduced per-  
19 token-memory compared to full Mamba fine-tuning, while achieving up to 81.5%  
20 of the ICL performance improvements (over zero-shot) of comparably fine-tuned  
21 Transformers.

## 22 1 Introduction

23 Innovating on previous state-space models (SSMs) [23, 11], Mamba [22] has been recently proposed  
24 as an accurate, sub-quadratic alternative to Transformer large language models (LLMs). Mamba was  
25 initially shown to greatly outperform comparable Transformer LLMs [5] across a large number of  
26 standard natural language benchmarks. Subsequently, pretrained Mamba models have been widely  
27 adapted across different data modalities [42, 65, 36, 46, 37], tasks [60, 62, 48, 63, 57, 37, 2], and  
28 architectures [1, 45, 40].

29 However, despite such rapid and widespread adaptation, evaluation of Mamba’s ability to perform  
30 standard downstream learning abilities exhibited by Transformer-based LLMs have either not been  
31 extensively conducted on *standard natural benchmarks* or are completely lacking. For instance, while  
32 recent works [45, 19, 30] have evaluated Mamba’s ability to perform in-context learning (ICL), such  
33 studies focused extensively on either non-natural tasks [30, 17] or non-standard benchmarks [25].

---

\*Equal Contribution

34 Furthermore, evaluation of Mamba’s mixed-precision fine-tuning (MPFT) and performance efficient  
35 fine-tuning (PEFT) capabilities are currently lacking. For the former, MPFT (and, by extension,  
36 mixed-precision inference) are made difficult due to potential sensitivities of Mamba’s recurrent  
37 dynamics, where [21, 29] suggest full precision (FP32) is required to perform stable training. For  
38 the latter, PEFT via standard low-rank adaptation (LoRA) [28] is made difficult within Mamba’s  
39 SSM layer (referred to herein as the MambaBlock) due highly customized SSM CUDA kernels which  
40 provide competitive performance to attention-based speedups [10] at the cost of standard adapter  
41 support. However, PEFT and MPFT are arguably two of the most widely utilized techniques for LLM  
42 alignment [53] and customization [55], and are typically combined to drastically decrease hardware  
43 demands needed to fine-tune modern LLMs [12].

44 Herein, we extensively explore Mamba’s downstream learning capabilities across standard natural  
45 benchmarks. For ICL, we show that, in contrast to recent non-standard studies showing Mamba  
46 models rival state-of-the-art (SOTA) LLMs of similar parameter counts, **the pretrained benefits of**  
47 **Mamba few-shot learning are significantly less than comparable Transformer LLMs across**  
48 **standard natural benchmarks**; averaged across the benchmarks and parameter counts in Table 1,  
49 **Mamba models only achieve 38% of the performance improvements (relative to zero-shot)**  
50 **of comparable Transformer models** from the Pythia suite [5]. However, we show in the sequel  
51 that **Mamba models can more than halve this gap through efficient fine-tuning**, achieving as  
52 much as 81.5% of the average few-shot learning improvement (relative to zero-shot) of comparable  
53 Transformers.

54 For MPFT, we leverage theory from dynamical systems to show that small input changes in a  
55 MambaBlock do not lead to exponentially deviating outputs. Empirically, we validate this theoretical  
56 result; compared to full-precision, deviations due to mixed-precision for Mamba inference and  
57 fine-tuning are on par with those demonstrated by Transformer LLMs (Section 6). For PEFT, we  
58 show that by targeting the largest memory buffer exploited by Mamba’s highly customized CUDA  
59 kernels, LoRA may be used for extremely efficient fine-tuning, while simultaneously regularizing  
60 the majority of Mamba’s SSM parameters via weight tying. We show that this leads to extremely  
61 efficient PEFT, resulting in up to 2.15 times faster training and 65.5% reduced memory compared to  
62 the largest evaluated Mamba model without MPFT or PEFT.

## 63 2 Background

64 **Downstream learning for LLMs.** Since the release of the Transformer architecture [54], attention-  
65 based LLMs have exhibited several downstream learning abilities—in particular, PEFT, MPFT, and  
66 ICL—which allow the rapid adaptation of foundation models towards specific applications. PEFT using  
67 adapters [24] allows a large pretrained model to be efficiently adapted for a particular downstream  
68 task by freezing the full model and training only a small number of extra parameters. Arguably the  
69 most widely used such PEFT method is LoRA [28], which injects trainable low-rank matrices into  
70 Transformer layers to approximate weight updates.

71 To further decrease the computational demands necessary for LLM fine-tuning and inference, MPFT  
72 via mixed-precision (i.e., FP16 or BF16) [31, 43] and quantized low-precision [12] have proven  
73 effective strategies to reduce GPU memory and runtime requirements without deleterious effects on  
74 downstream performance [12, 59]. Additionally, mixed-precision approaches have paved the way for  
75 hardware-aware optimizations within the self-attention module [10], greatly mitigating the quadratic  
76 complexity of Transformer LLMs. Together, PEFT and MPFT have created a rich ecosystem with  
77 which varying combinations of these approaches may be used to meet the computational constraints  
78 of a given training system. We note that post-fine-tuning quantization approaches [13] may be further  
79 used to decrease Transformer LLM computational demands, but such approaches are not considered  
80 in this work.

81 ICL provides an adaptable alternative to fine-tuning. Rather than fine-tune the LLM directly, ICL  
82 augments a prompt with  $n$  relevant examples (called *shots*) preceding the query of interest. Given  
83 sufficiently large models and pretraining data [8, 58], Transformer LLMs have proven adept at  
84 learning new concepts on the fly provided such few-shot prompting. However, it is worth noting  
85 that ICL inference time increases dramatically as the number of shots grows (due to self-attention’s  
86 quadratic complexity) and PEFT (when possible) is known to produce more accurate downstream  
87 learning results [8, 41].

Table 1: In-context learning performance for pretrained Mamba and Pythia models. Models are collected into parameter classes for head-to-head comparison using the groupings in [22]. Model checkpoints were evaluated on all benchmarks and few-shot settings using the LM evaluation harness from Eleuther AI [16]. LAMBADA zero-shot is more effective for the model sizes considered (further discussed in [61, 8]) and thus excluded from few-shot performance averages. Highlighted in bold is the top-performing few-shot learner per benchmark and model grouping.

Model	$N$ -shot	LAMBADA ppl ↓	LAMBADA acc ↑	HellaSwag acc ↑	PIQA acc ↑	Arc-E acc ↑	Arc-C acc ↑	WinoGrande acc ↑	0-shot incr. Mean % ↑
Mamba 130M	0	<b>16.07</b>	<b>44.3</b>	<b>35.3</b>	64.5	48.0	<b>24.2</b>	<b>44.8</b>	–
	1	19.34	38.3	35.2	64.3	47.1	23.5	51.3	-1.4
	3	23.13	35.4	35.1	<b>65.1</b>	<b>49.0</b>	24.0	50.7	-0.2
	5	24.38	36.2	34.8	64.9	49.2	23.9	50.5	-0.5
-----									
Pythia 160M	0	38.20	32.7	30.2	61.8	43.4	23.8	51.0	–
	1	47.21	28.2	30.6	62.2	43.4	23.7	49.3	-0.4
	3	63.70	24.7	30.5	61.9	44.8	22.9	51.3	<b>0.1</b>
	5	66.30	25.3	30.4	62.6	43.4	23.1	50.8	-0.2
-----									
Mamba 370M	0	<b>8.14</b>	<b>55.6</b>	<b>46.5</b>	69.5	55.0	27.9	55.5	–
	1	9.74	49.8	45.9	69.3	57.4	26.5	54.6	-0.8
	3	10.89	48.5	46.2	<b>69.6</b>	<b>58.7</b>	<b>28.5</b>	53.6	1.0
	5	11.36	48.5	46.2	69.4	58.3	28.0	<b>56.0</b>	1.3
-----									
Pythia 410M	0	10.83	51.5	40.6	66.9	52.0	24.1	53.4	–
	1	12.26	47.1	40.5	68.0	53.8	25.6	52.4	1.8
	3	14.39	43.2	40.9	67.9	55.1	26.9	54.0	<b>4.2</b>
	5	14.62	44.1	40.8	68.1	54.6	26.6	53.4	3.5
-----									
Mamba 790M	0	<b>6.01</b>	<b>61.7</b>	<b>55.1</b>	72.1	61.2	29.6	56.0	–
	1	7.06	56.2	54.5	<b>72.5</b>	63.3	30.1	56.9	1.4
	3	8.05	54.8	54.2	72.2	63.4	31.6	57.2	2.4
	5	8.83	53.4	54.6	<b>72.5</b>	<b>64.6</b>	<b>32.1</b>	<b>57.5</b>	<b>3.4</b>
-----									
Pythia 1B	0	7.92	56.3	47.2	70.7	57.0	27.0	53.4	–
	1	8.99	51.8	47.3	70.7	57.1	28.2	53.4	1.0
	3	10.48	48.2	47.5	71.2	59.2	28.0	54.3	2.2
	5	10.86	48.4	47.3	71.4	58.7	28.4	53.1	1.9
-----									
Mamba 1.4B	0	<b>5.04</b>	<b>65.0</b>	<b>59.1</b>	74.2	65.5	32.9	58.6	–
	1	5.83	60.6	58.20	<b>74.7</b>	64.5	33.0	61.2	-0.5
	3	6.62	58.9	58.8	73.7	66.1	34.4	60.9	0.6
	5	6.98	58.4	59.0	74.0	<b>66.4</b>	<b>35.5</b>	60.5	1.4
-----									
Pythia 1.4B	0	6.09	61.7	52.1	70.9	60.5	28.5	57.4	–
	1	6.96	56.3	52.1	71.4	62.0	29.5	57.5	1.4
	3	7.89	54.4	52.6	70.9	63.9	31.1	56.8	2.9
	5	8.02	54.4	52.8	71.0	63.2	31.3	57.8	<b>3.3</b>
-----									
Mamba 2.8B	0	<b>4.23</b>	<b>69.2</b>	<b>66.2</b>	75.2	69.7	36.3	63.4	–
	1	5.01	63.9	65.7	75.5	69.8	37.2	63.7	0.6
	3	5.53	63.0	65.5	75.2	70.8	38.1	<b>64.8</b>	1.6
	5	5.70	62.7	<b>66.2</b>	<b>76.2</b>	<b>70.9</b>	<b>38.3</b>	64.6	2.1
-----									
Pythia 2.8B	0	5.04	64.7	59.3	73.9	64.2	32.9	59.8	–
	1	5.66	60.9	59.4	73.8	66.8	34.8	59.0	1.7
	3	6.20	59.1	59.9	74.7	67.4	34.9	60.8	2.9
	5	6.52	59.1	60.2	74.5	67.1	35.0	61.3	<b>3.1</b>

88 **State-space Models.** Structured state-space sequence (S4) models [23, 14] are SSMs which leverage  
 89 linear time-invariant (LTI) systems to combine the computational advantages of Transformers—i.e.,  
 90 highly parallelizable training—and recurrent neural networks (RNNs)—i.e., subquadratic autoregressive  
 91 inference using recurrency. Within the S4 layer, an input signal is discretized and LTI parameters  
 92 representing the input’s latent dynamics are learned. Owing to the S4 block’s latent dynamics being  
 93 LTI, the S4 block’s output may be thus compactly represented as a single convolution between the  
 94 input and an *SSM convolution kernel* (a matrix whose entries are products of LTI learnable parameters  
 95 resulting from unrolling the state-space equations). However, despite hardware efficiency and  
 96 long-dependency-modeling improvements, LTI-based S4 models remained inferior to Transformers  
 97 of comparable parameter-sizes for natural language tasks, even when augmenting S4 layers with  
 98 attention-layers for hybrid architectures [22].

99 Innovating on these previous S4 approaches, Mamba utilizes time-varying parameters to model  
 100 latent dynamics, thus broadening the ability to capture nuanced changes evolving in discrete-time.  
 101 Without LTI dynamics, however, the input-output representation via the SSM convolution kernel is no  
 102 longer applicable, thus voiding previous hardware-aware S4 optimizations [14]. To enable hardware  
 103 efficiency with time-varying SSM parameters, [22] thus introduced extensively customized CUDA  
 104 kernels which implement highly parallelized prefix sums to compute recurrent states.

### 105 3 Mamba state-space models

106 For model dimension  $d$  and maximum input sequence length  $T$ , the MambaBlock defines state-space  
 107 parameters  $\mathbf{A}, \mathbf{B}_t, \mathbf{C}_t, \mathbf{\Delta}_t \in \mathbb{R}^{d \times d}$  for  $t \in \{1, \dots, T\}$ . The matrix  $\mathbf{\Delta}_t$  controls the discrete step-  
 108 size. Given an input sequence  $\mathbf{u}_1, \dots, \mathbf{u}_T \in \mathbb{R}^d$ , the following linear mapping through latent states  
 109  $\mathbf{x}_1, \dots, \mathbf{x}_T \in \mathbb{R}^d$  is used to produce the output  $\mathbf{y}_1, \dots, \mathbf{y}_T \in \mathbb{R}^d$ :

$$\mathbf{x}_t = \bar{\mathbf{A}}_t \mathbf{x}_{t-1} + \bar{\mathbf{B}}_t \mathbf{u}_t \quad (1)$$

$$\mathbf{y}_t = \bar{\mathbf{C}}_t \mathbf{x}_t, \quad (2)$$

110 where  $\bar{\mathbf{A}}_t = \text{softplus}(\text{Linear}(\mathbf{\Delta}_t)) \in \mathbb{R}^{d \times d}$ ,  $\bar{\mathbf{A}}_t = \exp(\bar{\mathbf{\Delta}}_t \mathbf{A})$  and  $\bar{\mathbf{B}}_t = \mathbf{A}^{-1}(\bar{\mathbf{A}} - \mathbf{I})\mathbf{B}_t$ . In  
 111 practice,  $\mathbf{A}, \mathbf{B}_t, \mathbf{C}_t$  and  $\mathbf{\Delta}_t$  are diagonal matrices.

112 **Hardware-aware optimizations.** As matrices  $\mathbf{B}_t, \mathbf{C}_t$  and  $\mathbf{\Delta}_t$  are time-varying, S4 optimizations via  
 113 the SSM convolution kernel [11] are no longer applicable. However, by diagonality, each dimension  
 114 may be computed in parallel. Furthermore, the recurrence along every dimension is a prefix sum (also  
 115 called a *scan*), which is highly parallelizable [7]. [15] thus capitalizes on this through extensively  
 116 customized CUDA kernels wherein the majority of temporal variables are carefully laid out in a large  
 117 buffer of GPU memory and manipulated. Instantiated as a PyTorch linear layer’s weight matrix, this  
 118 memory buffer  $\mathbf{W} \in \mathbb{R}^{n \times 3d}$  is used to store and access the diagonal elements of  $\mathbf{B}_t, \mathbf{C}_t$  and  $\mathbf{\Delta}_t$  for  
 119 all  $t \in \{1, \dots, T\}$ , such that

$$\mathbf{W}[t-1, : d] = \text{diag}(\mathbf{\Delta}_t), \mathbf{W}[t-1, d : 2d] = \text{diag}(\mathbf{B}_t), \mathbf{W}[t-1, 2d : 3d] = \text{diag}(\mathbf{C}_t), \quad (3)$$

120 where  $\mathbf{W}[0, : d] = \text{diag}(\mathbf{\Delta}_1)$ ,  $\mathbf{W}[n-1, d : 2d] = \text{diag}(\mathbf{B}_T)$ , and so on.

121 The customized Mamba prefix scan kernel heavily relies on this memory layout to optimize the  
 122 access pattern of  $\mathbf{W}$  in Equations 5 and 6. We note that, rather than adjusting Mamba’s low-level  
 123 CUDA kernels themselves to integrate LoRA within the highly optimized prefix scan, we can instead  
 124 directly target  $\mathbf{W}$ . Doing so, we have the following, where the proof is available in Appendix A.

125 **Theorem 1.** Consider the weight matrix  $\mathbf{W}$  of a MambaBlock from Equation 3. Targeting  $\mathbf{W}$  for  
 126 LoRA during fine-tuning ties adaptation weights across  $\mathbf{B}_t, \mathbf{C}_t$  and  $\mathbf{\Delta}_t$ .

### 127 4 Stable dynamics in the MambaBlock

128 The Mamba foundation models were pretrained in full FP32 precision. Consequently, official Mamba  
 129 implementations have cautioned against fine-tuning or training in reduced precision [21, 29], with  
 130 potential sensitivities of MambaBlock recurrent dynamics remaining an open question. We answer  
 131 the latter using theory from dynamical systems. For Mamba’s discrete dynamic system in Equations 5  
 132 and 6, define

$$\mathbf{x}_t = F_\theta(\mathbf{x}_{t-1}, \mathbf{u}_t), \quad (4)$$

133 where  $\theta$  denotes the time-varying parameters described in Section 3. For input sequence  $\mathbf{u}_1, \dots, \mathbf{u}_T$   
 134 and initial latent state vector  $\mathbf{x}_0$ , we thus write

$$\mathbf{x}_T = F_\theta(F_\theta(\dots F_\theta(\mathbf{x}_0, \mathbf{u}_1))) := F_\theta^{T-1}(\mathbf{x}_0, \mathbf{u}_1).$$

135 The rate of divergence between two scalar  $\varepsilon$ -close inputs to a discrete dynamical system is bounded  
 136 by the system’s maximal Lyapunov exponent  $\lambda_{\max}$  [44]. Given  $\lambda_{\max}$  and two initial values  $(\mathbf{x}_0, \mathbf{u}_1)$   
 137 and  $(\mathbf{x}_0 + \varepsilon, \mathbf{u}_1 + \varepsilon)$ , the maximum deviation between these points grows as [33, 50]:

$$\max |F_\theta^N(\mathbf{x}_0, \mathbf{u}_1) - F_\theta^N(\mathbf{x}_0 + \varepsilon, \mathbf{u}_1 + \varepsilon)| \in \mathcal{O}(\varepsilon \exp(N\lambda_{\max})).$$

138 Thus, when  $\lambda_{\max} > 0$ , nearby trajectories exponentially separate and, when  $\lambda_{\max} \leq 0$ , nearby  
 139 trajectories ultimately converge to the same fixed point or periodic cycles.

140 The maximal Lyapunov exponent is defined as

$$\lambda_{\max} := \lim_{T \rightarrow \infty} \frac{1}{T} \log \left\| \prod_{t=0}^{T-1} \frac{\partial \mathbf{x}_t}{\partial \mathbf{x}_{t-1}} \right\|_2,$$

141 where  $\|\cdot\|_2$  denotes the spectral norm for matrices. For an arbitrary MambaBlock, we prove the  
 142 following:

143 **Theorem 2.** *Let  $(\mathbf{x}_{t-1}, \mathbf{u}_t)$  be the latent state and input at an arbitrary time  $t \in \{1, \dots, T\}$  within a*  
 144 *MambaBlock. Then small changes  $(\mathbf{x}_{t-1} + \varepsilon, \mathbf{u}_t + \varepsilon)$  produce deviations which are exponentially non-*  
 145 *increasing over discrete-time. That is,  $\max |F_\theta^N(\mathbf{x}_{t-1}, \mathbf{u}_t) - F_\theta^N(\mathbf{x}_{t-1} + \varepsilon, \mathbf{u}_t + \varepsilon)| \in \mathcal{O}(\varepsilon \exp(N\zeta))$ ,*  
 146 *for some scalar  $\zeta \leq 0$ .*

147 The proof of Theorem 2 is available in Appendix B, where the maximal Lyapunov exponent for an  
 148 arbitrary MambaBlock is first proven to be non-positive. The main result subsequently follows.

149 **Consequences for automatic mixed-precision.** During a forward pass, automatic mixed-precision  
 150 (AMP) saves time and memory by computing forward activations in half-precision (FP16 or BF16).  
 151 During a backward pass, AMP computes gradients in half-precision and up-casts to full-precision  
 152 prior to updating. In contrast to full-precision fine-tuning, MPFT within the MambaBlock thus results  
 153 in small differences to the inputs  $\mathbf{u}_1, \dots, \mathbf{u}_T$  fed into the SSM scan (which are passed through a  
 154 SwiGLU),  $\hat{\Delta}_t$  (which is passed through a softplus), and the gradients calculated during training.

155 For a discrete dynamical system with  $\lambda_{\max} > 0$ , changes due to AMP compound after repeated  
 156 expansion of the recurrent state, thus leading to exponential deviations between quantities calculated  
 157 using mixed- versus full-precision. We note that Transformers are not recurrent, and thus not  
 158 susceptible to such issues. Yet, just as differences introduced by quantization/mixed-precision produce  
 159 output differences in Transformer results, differences are expected in Mamba results using different  
 160 precision strategies. However, by Theorem 2, such differences do not exponentially compound over  
 161 discrete-time within the MambaBlock.

## 162 5 Related Work

163 Several recent works [45, 19, 30, 40] have studied Mamba’s ability to perform ICL. However, none  
 164 of these have extensively studied Mamba’s ICL capabilities either on standard NLP benchmarks or on  
 165 pure MambaBlock foundation models. In particular, foundational Mamba models’ ICL abilities were  
 166 tested in [45] to learn simple function classes (e.g., logistic regression and decision trees [17]) and in  
 167 [19] to learn non-standard NLP benchmarks (i.e., task vectors [25]). While [45, 19] report Mamba’s  
 168 ICL abilities rival SOTA Transformers, their utilized benchmarks were proposed as supplemental  
 169 ICL studies after Transformer LLMs’ success on standard NLP benchmarks [8]. Indeed, direct  
 170 evaluation of Mamba foundation models on standard NLP benchmarks does not lead to higher gains  
 171 over zero-shot performance relative to comparable Transformer LLMs (demonstrated in Table 1).

172 Lyapunov exponents have previously been considered for classic RNN structures (e.g., vanilla  
 173 RNNs, LSTMs, GRUs, PLRNNs, etc.) [44, 56], to determine when such models exhibit chaotic  
 174 dynamics and the impact on the exploding/vanishing gradient phenomena\*. For more recent S4 neural

\*We note that this continues a long line of research exploring RNNs sensitivity to initial conditions and their subsequent ability to produce chaotic output [47, 34, 3, 4], although previous work did not leverage Lyapunov exponents.

175 models, [18] used Hurwitz matrices to characterize the numerical stability of linear time-invariant  
 176 (LTI) S4 models. However, such analysis is not applicable to time-varying models, such as Mamba,  
 177 nor does it characterize the effects of sensitive dependence on initial conditions (e.g., divergence of  
 178 two  $\varepsilon$  close inputs). To the best of our knowledge, no previous works have used Lyapunov exponents  
 179 to explore the effects of mixed-precision on recurrent neural models or Mamba architectures.

180 As in [22], the majority of subsequent Mamba works have focused on pretraining MambaBlocks using  
 181 full precision [65, 62, 1, 40]. Notably, the official implementation of Jamba [40], the Transformer-  
 182 Mamba hybrid, supports mixed- and 8-bit precision, but avoids MambaBlocks when applying such  
 183 quantization [32]. Similarly, the official Mamba sources advise using full precision within the  
 184 MambaBlock [29, 21], cautioning against using mixed-precision due to potential recurrent sensitivities.  
 185 To the best of our knowledge, no existing works have either theoretically explored the effects small  
 186 input changes (e.g., due to mixed-precision) have on Mamba’s recurrent dynamics, empirically  
 187 explored such effects downstream impact on fine-tuning and inference, or explored pure Mamba  
 188 networks fine-tuning abilities relative to Transformer LLMs.

## 189 6 Experiments

190 To demonstrate the implications of Theorem 2, we explore the performance difference between  
 191 running inference with full-precision pretrained weights and using mixed-precision (FP16 and BF16)  
 192 weights. **Model performance is measured as percent accuracy** using the MMLU [26] dataset.  
 193 The difference in model performance is reported as the mean *divergence* (i.e., absolute difference)  
 194 between the original full-precision and respective mixed-precision model, averaged over {0, 1, 3,  
 195 5}-shot percent accuracy. Thus, **a divergence greater than one denotes an average difference**  
 196 **greater than one entire percentage of accuracy.**

197 Mamba pretrained checkpoints are compared to pretrained Transformer models of similar parameter  
 198 counts and no more than  $\sim 300$ B total pretraining tokens (Pythia [5], OLMo [20] 336B-token  
 199 checkpoint, and Phi 1.5 [39]). We note that Pythia and Mamba models were both pretrained using  
 200 the same corpus [15], allowing the fairest comparison between SSMs and Transformers. To limit  
 201 extraneous numerical effects within experiments (e.g., due to parameter aggregation across multiple  
 202 GPUs), all models were run using a single GPU (Nvidia A10G, 24 GB total memory). All models  
 203 were evaluated using the LM evaluation harness from Eleuther AI [16]. Further experimental details  
 204 are available in Appendix C. The results are available in Table 2.

Table 2: Mean full-precision (FP32) divergence in MMLU performance for mixed-precision inference. Divergence is averaged over {0, 1, 3, 5}-shot performance. Pretrained checkpoints are used for Mamba (M), Pythia (P), OLMo [20], and Phi-1.5 [39] (Phi) models.

Model	M	P	M	P	M	P	OLMo	M	P	Phi	M	P
Size	130m	160m	370m	410m	790m	1b		1.4b	1.5b		2.8b	
FP16 $\mu$	0.03	0.35	0.05	0.06	0.21	0.05	0.04	0.04	0.07	0.03	0.15	0.12
BF16 $\mu$	0.05	1.45	0.20	0.20	0.66	0.16	0.13	0.31	0.13	1.05	1.17	0.11

205 From Table 2, inferring in Pythia using FP16 and BF16 result in an average 0.13 and 0.41 full-  
 206 precision divergence, respectively. Mamba displays similar averages in comparison: inferring in  
 207 Mamba using FP16 and BF16 result in an average 0.10 and 0.48 divergence, respectively. Interestingly,  
 208 both SSM and Transformer architectures exhibit *large divergence spikes*—i.e., mean divergence greater  
 209 than a percentage point—when using BF16, which occurs once for Mamba and Phi 1.5 models and  
 210 twice for Pythia models. In the following, we show that such spikes may be mitigated for Mamba  
 211 SSMs by combining mixed-precision with parameter-efficient adapters during fine-tuning.

212 **Non-divergent Mamba fine-tuning.** We next explore the implications of Theorem 2 on fine-tuning,  
 213 wherein mixed-precision is especially critical; MPFT combined with PEFT adapters have been shown  
 214 to drastically reduce Transformer fine-tuning times [12]. We are thus interested in the divergence  
 215 between Mamba models fully fine-tuned (i.e., no adapters, all model weights are trained) in full-  
 216 precision and models fine-tuned using mixed-precision and/or PEFT adapters. We focus on utilizing  
 217 LoRA [28], which is arguably the most widely used PEFT framework for LLMs.

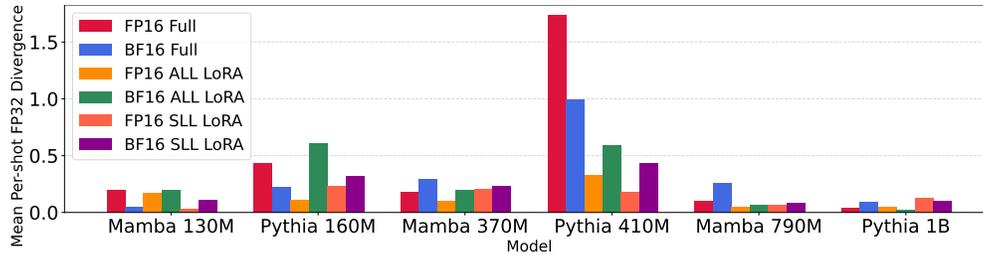


Figure 1: Mean full-precision (FP32) divergence in MMLU performance for Mamba and Pythia models. Models are fine-tuned over the Alpaca dataset [51] using different combinations of MPFT and PEFT. Full fine-tuning (i.e., no PEFT adapters) is denoted as Full.

218 Using the Alpaca dataset [51], Mamba 160M, 410M, and 790M models are fine-tuned for three epochs  
 219 with a maximum sequence length of 512. We denote the targeting of all linear layers (ALL) for LoRA  
 220 as *ALL LoRA*, the targeting of a subset of linear layers (SLL) for LoRA as *SLL LoRA*, and no adapters  
 221 as *Full* (i.e., full fine-tuning). Both ALL and SLL LoRA adapt the large memory buffer described in  
 222 Theorem 1.

223 Each fine-tuning run occurred on a single A10G GPU. To further limit extraneous numerical effects,  
 224 the same batch size is used for all FP32, FP16, and BF16 experiments for a given model size. While  
 225 this leads to hardware underutilization (i.e., non-saturated GPU memory for mixed-precision and  
 226 LoRA experiments), this is necessary to guarantee no divergence is due to differences in parameter  
 227 update schedules. For comparison, Pythia 160M, 410M, and 1B models are fine-tuned using the  
 228 same experimental setup. The training recipe for all models was adapted from [53], with the  
 229 AdamW\_torch optimizer and a cosine annealing schedule. Further experimental details are  
 230 available in Appendix C.

231 For each Mamba and Pythia model, Figure 1 shows the mean divergence calculated between the  
 232 respective FP32 Full and mixed-precision ALL/SLL LoRA fine-tuned models, averaged over {0, 1, 3,  
 233 5}-shot MMLU accuracy. Across mixed-precisions and adapter settings, Mamba displays comparable  
 234 divergences to Pythia models. E.g., for FP16, **Mamba demonstrates an average divergence of 0.1,**  
 235 **compared to 0.14 for Pythia.** Similarly, for BF16, **Mamba demonstrates an average divergence**  
 236 **of 0.18, compared to 0.28 for Pythia.** Importantly, Mamba models do not exhibit large deviation  
 237 spikes after fine-tuning (in contrast to Pythia models).

238 **Hardware throughput and memory-utilization improvements.** With comparable divergences  
 239 to Transformers and stable dynamics, we show that MPFT and PEFT may be used to significantly  
 240 increase GPU-training throughput for Mamba SSMS. To demonstrate such improvements, we utilize  
 241 the previous fine-tuning settings for the Alpaca dataset. However, we now adjust the batch size to  
 242 maximize throughput per MPFT and PEFT configuration.

243 For each MPFT and PEFT configuration, the *average tokens-per-second* (ATPS) is calculated as the  
 244 total tokens used for fine-tuning divided by total training time, and the *maximum memory-per-token*  
 245 (MMPT) is calculated as the maximum GPU memory utilization incurred (over the entire fine-tuning  
 246 run) divided by the total number of tokens in each mini-batch. Results are plotted in Figure 6.

247 Both throughput and memory utilization improve as the number of Mamba parameters increases  
 248 in Figure 6. **Compared to the full-precision full fine-tuning of Mamba 790M** (the largest model  
 249 supported by an A10G’s memory capacity), **evaluated MPFT and PEFT combinations result in**  
 250 **an average 2.15 times more training tokens-per-second while reducing per-token memory**  
 251 **utilization by an average 62.7%.** Across all model sizes, evaluated MPFT and PEFT combinations  
 252 result in an average 1.74 times more training tokens-per-second while reducing per-token memory  
 253 utilization by an average 47.2% compared to respective full-precision fine-tuned runs.

## 254 6.1 Fine-tuning narrows the ICL gap between Mamba and Transformers

255 We next explore how MPFT and PEFT affect Mamba ICL performance. All Mamba pretrained  
 256 models are instruction fine-tuned using ALL LoRA and the OpenHermes dataset [52] (which consists  
 257 of 242,000 supervised samples). We use the training recipe of [53], which includes BF16 utilization.

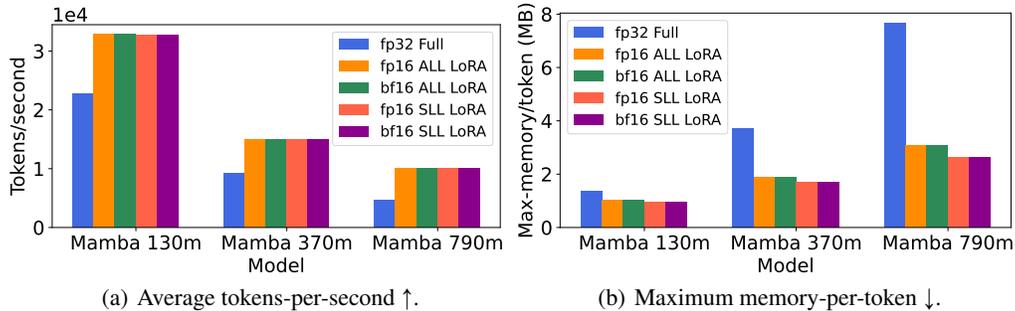


Figure 2: Timing and memory usage calculated Mamba model-sizes and PEFT combinations. Each model was trained using the Alpaca dataset [51] dataset for three epochs and maximum sequence length 512. For each PEFT combination, the batch size was tuned to maximize GPU occupancy.

258 Performance is evaluated using the datasets from Table 1–HellaSwag [64], PIQA [6], Arc-E [9],  
 259 Arc-C [9], and WinoGrande [49]–and report the *average improvement percentage* of {1, 3, 5}-shot  
 260 versus 0-shot (AIPSS). For comparison, Pythia pretrained models are instruction fine-tuned using the  
 same training recipe and ALL LoRA (i.e., all Pythia linear layers are adapted).

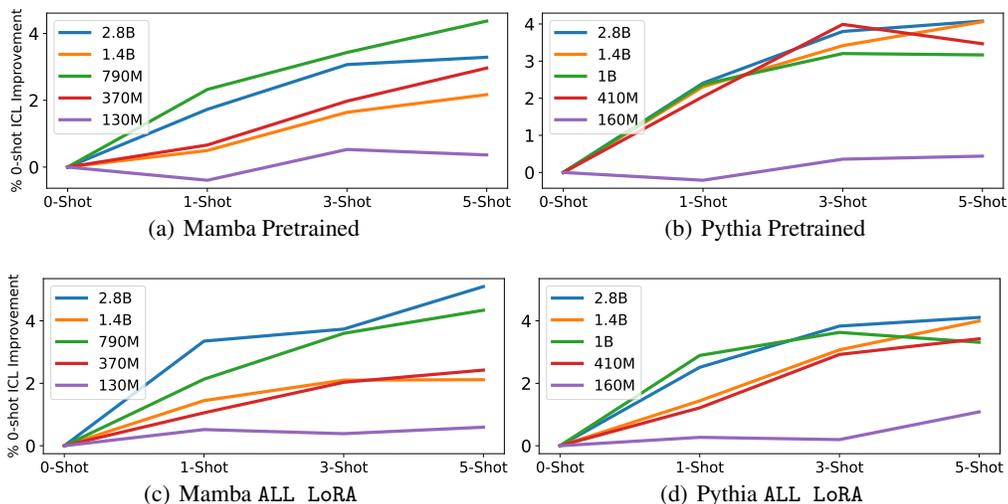


Figure 3: Fine-tuning narrows the ICL gap between Mamba and Pythia. ALL LoRA models were instruction fine-tuned on the OpenHermes [52] dataset for one epoch. Performance is reported as the average improvement percentage of {1, 3, 5}-shot versus 0-shot over five standard benchmarks.

261

262 Figure 3 displays AIPSS for pretrained and instruction fine-tuned Mamba and Pythia models. As  
 263 previously noted, pretrained Mamba models do not display similar ICL ability as comparable Pythia  
 264 models on the evaluated standard NLP benchmarks. In particular, Mamba 2.8B, the largest pretrained  
 265 Mamba model, displays inconsistent zero-shot improvements as the number of shots increase.  
 266 However, after fine-tuning, all Mamba models larger than Mamba 130M consistently improve in ICL  
 267 performance as the number of shots increase. Compared to Mamba pretrained models, which are only  
 268 capable of 38% of the AIPSS compared to similar pretrained Pythia models, fine-tuned ALL LoRA  
 269 Mamba models are capable of 81.5% of the AIPSS compared to similarly fine-tuned Pythia models.

270 **Fine-tuning robustness.** We show that Mamba is robust to the choice of PEFT hyperparameters. We  
 271 conduct an extensive hyperparameter search across the learning rate, LoRA dimension, and number of  
 272 warmup steps. From the Cartesian-product of these three parameters, 150 hyperparameter configura-  
 273 tions were sampled and used to fine-tune Mamba 370M over the Openhermes dataset. For comparison,  
 274 Pythia 410M is similarly fine-tuned using the same set of 150 hyperparameter configurations.

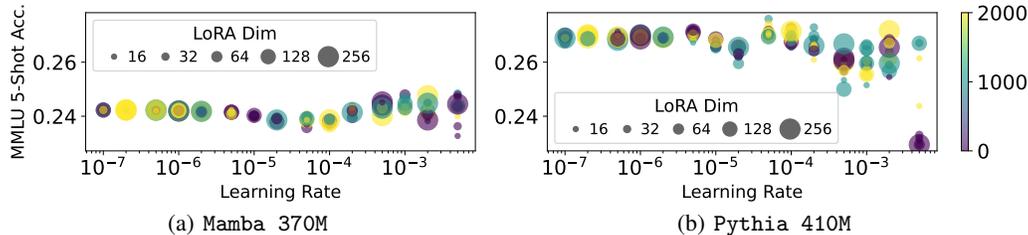


Figure 4: Fine-tuning hyperparameter search for OpenHermes. Each point is a different hyperparameter configuration. SLL LoRA was used for both models. The  $x$ -axis is the learning rate, the  $y$ -axis is resulting MMLU 5-shot performance, bubble size is the LoRA dimension, and the color is the number of warmup steps  $\in \{0, 1k, 2k\}$ .

275 The MMLU 5-shot performance for each of the 150 Mamba and Pythia fine-tuned models is displayed  
 276 in 6.1. Pythia 410M is capable of higher performance than Mamba 370M, where the average accuracy  
 277 for the former and the latter are 26.5% and 24.8%, respectively. However, Mamba 370M is much more  
 278 robust to the choice of hyperparameters, with a difference of 1.5% between the minimum (23.3%)  
 279 and maximum (24.8%). In contrast, Pythia 410M fine-tuned models display a large performance  
 280 difference of 4.7% between the minimum (22.9%) and maximum (27.6%).

## 281 7 Discussion

282 We’ve extensively explored Mamba’s downstream learning capabilities. Using dynamical systems  
 283 theory, we’ve shown that Mamba’s recurrent dynamics are robust to small input perturbations (contrary  
 284 to the current understanding of Mamba’s recurrent sensitivities). We’ve extensively confirmed this  
 285 result, showing that: a) Mamba inference is robust to changes due to mixed-precision, (b) Mamba  
 286 inference differences due to mixed-precision align with Transformers, (c) Mamba fine-tuning is robust  
 287 to changes due to mixed-precision and PEFT, and (d) differences in downstream performance for  
 288 Mamba due to MPFT and PEFT can be more robust than Transformers. Using both MPFT and PEFT,  
 289 we’ve shown that instruction fine-tuning Mamba SSMs greatly narrows the previously observed ICL  
 290 gap, going from only 38% (post pretraining) up to 81.5% (post fine-tuning) of the ICL abilities of  
 291 similar Transformers. Furthermore, we’ve shown that combining MPFT and PEFT can more than  
 292 halve training time and nearly triple memory efficiency for Mamba models.

293 There are significant avenues for future work. In particular, adapting Mamba’s CUDA kernels to  
 294 support more aggressive low-precision PEFT methods [12] would further decrease the hardware  
 295 needed to train Mamba models, while providing additional speedups. Furthermore, while the largest  
 296 pure Mamba model contains 2.8B parameters, the training speedups and improved memory utilization  
 297 described herein may be applied to more efficiently pretrain larger pure Mamba SSMs (e.g., 7B  
 298 parameters and greater), where Mamba models may better manifest emergent abilities previously  
 299 displayed by Transformers (or even manifest previously unobserved abilities).

300 **Limitations.** While we explored the use of LoRA for Mamba models, many other PEFT adapters  
 301 exist [41, 38, 27, 35]. Furthermore, while mixed-precision using FP16 and BF16 were explored,  
 302 lower-precision methods exist [12] (which may be enabled by adapting Mamba’s highly customized  
 303 CUDA kernels). Both are interesting directions for future work. Finally, our timing and memory  
 304 usage experiments using Alpaca did not consider the largest two Mamba models (1.4B and 2.8B) due  
 305 to their exceeding A10G memory capacity for FP32 full fine-tuning.

306 **Broader Impact.** The Mamba models considered are all LLMs, and thus have the same potential  
 307 positive and negative societal impacts as other LLMs (e.g., hallucinations). Furthermore, fine-tuning  
 308 is known to possibly erode existing LLM guardrails, and thus our methods may be adapted for this  
 309 fine-tuning use case (as is the case for all PEFT and MPFT methods). However, our work improves the  
 310 quality of Mamba models for downstream applications, which may be adapted for all positive LLM  
 311 applications in society (e.g., personal assistants, task automation, code completion, etc.). Finally, our  
 312 work decreases the computational constraints required to train and inference Mamba SSMs, which  
 313 has implications for green ML (e.g., decreased CO2 emissions, positive climate change impact, etc.).  
 314 410 GPU days were used to produce the results for this paper.

315 **References**

- 316 [1] Quentin Anthony, Yury Tokpanov, Paolo Glorioso, and Beren Millidge. Blackmamba: Mixture of experts  
317 for state-space models. *arXiv preprint arXiv:2402.01771*, 2024.
- 318 [2] Ali Behrouz and Farnoosh Hashemi. Graph mamba: Towards learning on graphs with state space models.  
319 *arXiv preprint arXiv:2402.08678*, 2024.
- 320 [3] Nils Bertschinger and Thomas Natschläger. Real-time computation at the edge of chaos in recurrent neural  
321 networks. *Neural computation*, 16(7):1413–1436, 2004.
- 322 [4] Nils Bertschinger, Thomas Natschläger, and Robert Legenstein. At the edge of chaos: Real-time computa-  
323 tions and self-organized criticality in recurrent neural networks. *Advances in neural information processing*  
324 *systems*, 17, 2004.
- 325 [5] Stella Biderman, Hailey Schoelkopf, Quentin Gregory Anthony, Herbie Bradley, Kyle O’Brien, Eric  
326 Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward Raff, et al. Pythia:  
327 A suite for analyzing large language models across training and scaling. In *International Conference on*  
328 *Machine Learning (ICML)*, pages 2397–2430. PMLR, 2023.
- 329 [6] Yonatan Bisk, Rowan Zellers, Jianfeng Gao, Yejin Choi, et al. Piqa: Reasoning about physical common-  
330 sense in natural language. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34,  
331 pages 7432–7439, 2020.
- 332 [7] Guy E Blelloch. Prefix sums and their applications. 1990.
- 333 [8] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind  
334 Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners.  
335 *Advances in neural information processing systems*, 33:1877–1901, 2020.
- 336 [9] Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind  
337 Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge. *arXiv preprint*  
338 *arXiv:1803.05457*, 2018.
- 339 [10] Tri Dao, Dan Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. Flashattention: Fast and memory-efficient  
340 exact attention with io-awareness. *Advances in Neural Information Processing Systems*, 35:16344–16359,  
341 2022.
- 342 [11] Tri Dao, Daniel Y Fu, Khaled K Saab, Armin W Thomas, Atri Rudra, and Christopher Ré. Hungry hungry  
343 hippos: Towards language modeling with state space models. In *Proceedings of the 11th International*  
344 *Conference on Learning Representations (ICLR)*, 2023.
- 345 [12] Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. Qlora: Efficient finetuning of  
346 quantized llms. *Advances in Neural Information Processing Systems*, 36, 2024.
- 347 [13] Elias Frantar, Saleh Ashkboos, Torsten Hoeffer, and Dan Alistarh. Gptq: Accurate post-training quantization  
348 for generative pre-trained transformers. 2023.
- 349 [14] Daniel Y Fu, Tri Dao, Khaled Kamal Saab, Armin W Thomas, Atri Rudra, and Christopher Re. Hungry  
350 hungry hippos: Towards language modeling with state space models. In *International Conference on*  
351 *Learning Representations (ICLR)*, 2023.
- 352 [15] Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang,  
353 Horace He, Anish Thite, Noa Nabeshima, et al. The pile: An 800gb dataset of diverse text for language  
354 modeling. *arXiv preprint arXiv:2101.00027*, 2020.
- 355 [16] Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster,  
356 Laurence Golding, Jeffrey Hsu, Alain Le Noac’h, Haonan Li, Kyle McDonell, Niklas Muennighoff, Chris  
357 Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron, Lintang Sutawika, Eric Tang,  
358 Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. A framework for few-shot language model evaluation,  
359 12 2023.
- 360 [17] Shivam Garg, Dimitris Tsipras, Percy S Liang, and Gregory Valiant. What can transformers learn in-  
361 context? a case study of simple function classes. *Advances in Neural Information Processing Systems*  
362 *(NeurIPS)*, 35:30583–30598, 2022.
- 363 [18] Karan Goel, Albert Gu, Chris Donahue, and Christopher Ré. It’s raw! audio generation with state-space  
364 models. In *International Conference on Machine Learning*, pages 7616–7633. PMLR, 2022.

- 365 [19] Riccardo Grazi, Julien Siems, Simon Schrodi, Thomas Brox, and Frank Hutter. Is mamba capable of  
366 in-context learning? *arXiv preprint arXiv:2402.03170*, 2024.
- 367 [20] Dirk Groeneveld, Iz Beltagy, Pete Walsh, Akshita Bhagia, Rodney Kinney, Oyvind Tafjord, Ananya Harsh  
368 Jha, Hamish Ivison, Ian Magnusson, Yizhong Wang, et al. Olmo: Accelerating the science of language  
369 models. *arXiv preprint arXiv:2402.00838*, 2024.
- 370 [21] Albert Gu and Tri Dao. *Mamba Precision Guidance*. "[https://github.com/state-spaces/mamba#](https://github.com/state-spaces/mamba#precision)  
371 [precision](https://github.com/state-spaces/mamba#precision)", 2023. "Accessed: 2024-04-25".
- 372 [22] Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. *arXiv preprint*  
373 *arXiv:2312.00752*, 2023.
- 374 [23] Albert Gu, Karan Goel, and Christopher Re. Efficiently modeling long sequences with structured state  
375 spaces. In *International Conference on Learning Representations (ICLR)*, 2022.
- 376 [24] Junxian He, Chunting Zhou, Xuezhe Ma, Taylor Berg-Kirkpatrick, and Graham Neubig. Towards a unified  
377 view of parameter-efficient transfer learning. In *International Conference on Learning Representations*  
378 *(ICLR)*, 2021.
- 379 [25] Roei Hendel, Mor Geva, and Amir Globerson. In-context learning creates task vectors. In *The 2023*  
380 *Conference on Empirical Methods in Natural Language Processing*, 2023.
- 381 [26] Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob  
382 Steinhardt. Measuring massive multitask language understanding. In *International Conference on Learning*  
383 *Representations*, 2020.
- 384 [27] Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea Ges-  
385 mundo, Mona Attariyan, and Sylvain Gelly. Parameter-efficient transfer learning for nlp. In *International*  
386 *conference on machine learning*, pages 2790–2799. PMLR, 2019.
- 387 [28] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and  
388 Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*,  
389 2021.
- 390 [29] Huggingface. *Mamba PEFT*. "[https://huggingface.co/docs/transformers/en/model\\_doc/](https://huggingface.co/docs/transformers/en/model_doc/mamba#peft-finetuning)  
391 [mamba#peft-finetuning](https://huggingface.co/docs/transformers/en/model_doc/mamba#peft-finetuning)", 2024. "Accessed: 2024-04-25".
- 392 [30] Samy Jelassi, David Brandfonbrener, Sham M Kakade, and Eran Malach. Repeat after me: Transformers  
393 are better than state space models at copying. *arXiv preprint arXiv:2402.01032*, 2024.
- 394 [31] Dhiraj Kalamkar, Dheevatsa Mudigere, Naveen Mellempudi, Dipankar Das, Kunal Banerjee, Sasikanth  
395 Avancha, Dharma Teja Vooturi, Nataraj Jammalamadaka, Jianyu Huang, Hector Yuen, et al. A study of  
396 bfloat16 for deep learning training. *arXiv preprint arXiv:1905.12322*, 2019.
- 397 [32] AI 21 Labs. *Jamba PEFT*. "<https://huggingface.co/ai21labs/Jamba-v0.1>", 2024. "Accessed:  
398 2024-04-25".
- 399 [33] Tanguy Laffargue, Khanh-Dang Nguyen Thu Lam, Jorge Kurchan, and Julien Tailleur. Large deviations of  
400 lyapunov exponents. *Journal of Physics A: Mathematical and Theoretical*, 46(25):254002, 2013.
- 401 [34] Thomas Laurent and James von Brecht. A recurrent neural network without chaos. In *Proceedings of the*  
402 *11th International Conference on Learning Representations (ICLR)*, 2017.
- 403 [35] Brian Lester, Rami Al-Rfou, and Noah Constant. The power of scale for parameter-efficient prompt tuning.  
404 *arXiv preprint arXiv:2104.08691*, 2021.
- 405 [36] Kai Li and Guo Chen. Spmamba: State-space model is all you need in speech separation. *arXiv preprint*  
406 *arXiv:2404.02063*, 2024.
- 407 [37] Lincan Li, Hanchen Wang, Wenjie Zhang, and Adelle Coster. Stg-mamba: Spatial-temporal graph learning  
408 via selective state space model. *arXiv preprint arXiv:2403.12418*, 2024.
- 409 [38] Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation. *arXiv*  
410 *preprint arXiv:2101.00190*, 2021.
- 411 [39] Yuanzhi Li, Sébastien Bubeck, Ronen Eldan, Allie Del Giorno, Suriya Gunasekar, and Yin Tat Lee.  
412 Textbooks are all you need ii: phi-1.5 technical report. *arXiv preprint arXiv:2309.05463*, 2023.

- 413 [40] Opher Lieber, Barak Lenz, Hofit Bata, Gal Cohen, Jhonathan Osin, Itay Dalmedigos, Erez Safahi, Shaked  
414 Meirum, Yonatan Belinkov, Shai Shalev-Shwartz, et al. Jamba: A hybrid transformer-mamba language  
415 model. *arXiv preprint arXiv:2403.19887*, 2024.
- 416 [41] Haokun Liu, Derek Tam, Mohammed Muqeeth, Jay Mohta, Tenghao Huang, Mohit Bansal, and Colin A  
417 Raffel. Few-shot parameter-efficient fine-tuning is better and cheaper than in-context learning. *Advances  
418 in Neural Information Processing Systems*, 35:1950–1965, 2022.
- 419 [42] Yue Liu, Yunjie Tian, Yuzhong Zhao, Hongtian Yu, Lingxi Xie, Yaowei Wang, Qixiang Ye, and Yunfan  
420 Liu. Vmamba: Visual state space model. *arXiv preprint arXiv:2401.10166*, 2024.
- 421 [43] Paulius Micikevicius, Sharan Narang, Jonah Alben, Gregory Diamos, Erich Elsen, David Garcia, Boris  
422 Ginsburg, Michael Houston, Oleksii Kuchaiev, Ganesh Venkatesh, et al. Mixed precision training. In  
423 *International Conference on Learning Representations (ICLR)*, 2018.
- 424 [44] Jonas Mikhaeil, Zahra Monfared, and Daniel Durstewitz. On the difficulty of learning chaotic dynamics  
425 with rnns. *Advances in Neural Information Processing Systems*, 35:11297–11312, 2022.
- 426 [45] Jongho Park, Jaeseung Park, Zheyang Xiong, Nayoung Lee, Jaewoong Cho, Samet Oymak, Kangwook  
427 Lee, and Dimitris Papailiopoulos. Can mamba learn how to learn? a comparative study on in-context  
428 learning tasks. *International Conference on Machine Learning (ICML)*, 2024.
- 429 [46] Changsheng Quan and Xiaofei Li. Multichannel long-term streaming neural speech enhancement for static  
430 and moving speakers. *arXiv preprint arXiv:2403.07675*, 2024.
- 431 [47] Antônio H Ribeiro, Koen Tiels, Luis A Aguirre, and Thomas Schön. Beyond exploding and vanishing  
432 gradients: analysing rnn training using attractors and smoothness. In *International conference on artificial  
433 intelligence and statistics (AISTATS)*, pages 2370–2380. PMLR, 2020.
- 434 [48] Jiacheng Ruan and Suncheng Xiang. Vm-unet: Vision mamba unet for medical image segmentation. *arXiv  
435 preprint arXiv:2402.02491*, 2024.
- 436 [49] Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Winogrande: An adversarial  
437 winograd schema challenge at scale. *Communications of the ACM*, 64(9):99–106, 2021.
- 438 [50] Hiroki Sayama. *Introduction to the modeling and analysis of complex systems*. Open SUNY Textbooks,  
439 2015.
- 440 [51] Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang,  
441 and Tatsunori B. Hashimoto. Stanford alpaca: An instruction-following llama model. [https://github.  
442 com/tatsu-lab/stanford\\_alpaca](https://github.com/tatsu-lab/stanford_alpaca), 2023.
- 443 [52] Teknium. *Openhermes*. "<https://huggingface.co/datasets/teknium/openhermes>", 2024. "Ac-  
444 cessed: 2024-04-25".
- 445 [53] Lewis Tunstall, Edward Beeching, Nathan Lambert, Nazneen Rajani, Kashif Rasul, Younes Belkada,  
446 Shengyi Huang, Leandro von Werra, Clémentine Fourrier, Nathan Habib, et al. Zephyr: Direct distillation  
447 of lm alignment. *arXiv preprint arXiv:2310.16944*, 2023.
- 448 [54] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz  
449 Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*,  
450 30, 2017.
- 451 [55] Kushala VM, Harikrishna Warriar, Yogesh Gupta, et al. Fine tuning llm for enterprise: Practical guidelines  
452 and recommendations. *arXiv preprint arXiv:2404.10779*, 2024.
- 453 [56] Ryan Vogt, Maximilian Puelma Touzel, Eli Shlizerman, and Guillaume Lajoie. On lyapunov exponents  
454 for rnns: Understanding information propagation using dynamical systems tools. *Frontiers in Applied  
455 Mathematics and Statistics*, 8:818799, 2022.
- 456 [57] Chloe Wang, Oleksii Tsepa, Jun Ma, and Bo Wang. Graph-mamba: Towards long-range graph sequence  
457 modeling with selective state spaces. *arXiv preprint arXiv:2402.00789*, 2024.
- 458 [58] Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama,  
459 Maarten Bosma, Denny Zhou, Donald Metzler, et al. Emergent abilities of large language models.  
460 *Transactions on Machine Learning Research*, 2022.
- 461 [59] Hao Wu, Patrick Judd, Xiaojie Zhang, Mikhail Isaev, and Paulius Micikevicius. Integer quantization for  
462 deep learning inference: Principles and empirical evaluation. *arXiv preprint arXiv:2004.09602*, 2020.

- 463 [60] Jianhao Xie, Ruofan Liao, Ziang Zhang, Sida Yi, Yuesheng Zhu, and Guibo Luo. Prommamba: Prompt-  
464 mamba for polyp segmentation. *arXiv preprint arXiv:2403.13660*, 2024.
- 465 [61] Sang Michael Xie, Aditi Raghunathan, Percy Liang, and Tengyu Ma. An explanation of in-context learning  
466 as implicit bayesian inference. In *International Conference on Learning Representations (ICLR)*, 2021.
- 467 [62] Zhaohu Xing, Tian Ye, Yijun Yang, Guang Liu, and Lei Zhu. Segmamba: Long-range sequential modeling  
468 mamba for 3d medical image segmentation. *arXiv preprint arXiv:2401.13560*, 2024.
- 469 [63] Yijun Yang, Zhaohu Xing, and Lei Zhu. Vivim: a video vision mamba for medical video object segmenta-  
470 tion. *arXiv preprint arXiv:2401.14168*, 2024.
- 471 [64] Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a machine really  
472 finish your sentence? In *Proceedings of the 57th Annual Meeting of the Association for Computational*  
473 *Linguistics*, pages 4791–4800, 2019.
- 474 [65] Lianghui Zhu, Bencheng Liao, Qian Zhang, Xinlong Wang, Wenyu Liu, and Xinggong Wang. Vision  
475 mamba: Efficient visual representation learning with bidirectional state space model. *International*  
476 *Conference on Machine Learning (ICML)*, 2024.

## 477 A Proof of weight-tying using LoRA in the MambaBlock

478 Due to the low-level nature of Mamba’s prefix scan optimizations (discussed in Section 3), standard  
479 use of LoRA adapters is made difficult within Mamba’s SSM-layer. E.g., while  $B_t, C_t$  and  $\Delta_t$  are  
480 conceptually PyTorch linear layers, their bundling in a contiguous memory block and careful manip-  
481 ulation makes appending a LoRA adapter on any of these individual matrices non-trivial (particularly,  
482 while respecting the highly specialized layout of each LoRA adapters targeted layer). However, we  
483 note that the overall design of the MambaBlock’s hardware optimizations may be leveraged to both  
484 efficiently learn the parameter-space for the majority of time-varying parameters (thus achieving  
485 PEFT) and regularize parameters during training (thus improving fine-tuning generalization).

486 **Theorem 1.** *Consider the weight matrix  $\mathbf{W}$  of a MambaBlock from Equation 3. Targeting  $\mathbf{W}$  for*  
487 *LoRA during fine-tuning ties adaptation weights across  $\mathbf{B}_t, \mathbf{C}_t$  and  $\Delta_t$ .*

488 *Proof.* Let  $r$  be the specified LoRA dimension. Targeting this matrix for LoRA results in the adapter

$$\begin{aligned}\tilde{\mathbf{W}} &= \mathbf{W} + \mathbf{W}' \\ &= \mathbf{W} + \mathbf{U}\mathbf{V},\end{aligned}$$

489 where  $\mathbf{U} \in \mathbb{R}^{n \times r}$ ,  $\mathbf{V} \in \mathbb{R}^{r \times 3d}$ , and  $\mathbf{W}$  is frozen during fine-tuning. Thus, for index  $[i, j]$ ,

$$\mathbf{W}'[i, j] = \sum_{k=0}^{r-1} \mathbf{U}[i, k] \mathbf{V}[k, j].$$

490 Recall the form of  $\mathbf{W}$ :

$$\mathbf{W}[t-1, : d] = \text{diag}(\Delta_t), \mathbf{W}[t-1, d : 2d] = \text{diag}(\mathbf{B}_t), \mathbf{W}[t-1, 2d : 3d] = \text{diag}(\mathbf{C}_t),$$

491 where  $\mathbf{W}[0, : d] = \text{diag}(\Delta_1)$ ,  $\mathbf{W}[n-1, d : 2d] = \text{diag}(\mathbf{B}_T)$ , and so on. For index  $[t-1, j]$ , we  
492 thus have

$$\begin{aligned}\tilde{\mathbf{W}}[t-1, j] &= \mathbf{W}[t-1, j] + \mathbf{W}'[t-1, j] \\ &= \mathbf{W}[t-1, j] + \sum_{k=0}^{r-1} \mathbf{U}[t-1, k] \mathbf{V}[k, j].\end{aligned}$$

493 Thus, the weights  $\mathbf{U}[t-1, :]$  are tied for any parameter  $\tilde{\mathbf{W}}[t-1, j], j \in \{1, \dots, 3d\}$ , which are used  
494 to adapt parameters  $\Delta_1, \mathbf{B}_t$ , and  $\mathbf{C}_t$ .

495 □

496 **B Mamba stable dynamics proof**

497 Recall the state-space parameters and equations for the MambaBlock;  $\mathbf{A}, \mathbf{B}_t, \mathbf{C}_t, \mathbf{\Delta}_t \in \mathbb{R}^{d \times d}$  for  
 498  $t \in \{1, \dots, n\} = [n]$ . Given an input sequence  $\mathbf{u}_1, \dots, \mathbf{u}_n \in \mathbb{R}^d$ , the following linear mapping  
 499 through latent states  $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$  is used to produce the output  $\mathbf{y}_1, \dots, \mathbf{y}_n \in \mathbb{R}^d$ :

$$\mathbf{x}_t = \bar{\mathbf{A}}_t \mathbf{x}_{t-1} + \bar{\mathbf{B}}_t \mathbf{u}_t \quad (5)$$

$$\mathbf{y}_t = \bar{\mathbf{C}}_t \mathbf{x}_t, \quad (6)$$

500 where  $\bar{\mathbf{\Delta}}_t = \text{softplus}(\text{Linear}(\mathbf{\Delta}_t)) \in \mathbb{R}^{d \times d}$ ,  $\bar{\mathbf{A}}_t = \exp(\bar{\mathbf{\Delta}}_t \mathbf{A})$ ,  $\bar{\mathbf{B}}_t = \mathbf{A}^{-1}(\bar{\mathbf{A}} - \mathbf{I})\mathbf{B}_t$ , and is the  
 501 set of non-negative real numbers. In practice,  $\mathbf{A}, \mathbf{B}_t, \mathbf{C}_t$  and  $\mathbf{\Delta}_t$  are diagonal matrices.

502 Furthermore, recall the following definitions:

$$\mathbf{x}_t = F_\theta(\mathbf{x}_{t-1}, \mathbf{u}_t)$$

503 where  $\theta$  denotes the aforementioned time-varying parameters. For input sequence  $\mathbf{u}_1, \dots, \mathbf{u}_T$  and  
 504 initial latent state value  $\mathbf{x}_0$ , we thus write

$$\mathbf{x}_T = F_\theta(F_\theta(\dots F_\theta(\mathbf{x}_0, \mathbf{u}_1))) := F_\theta^{T-1}(\mathbf{x}_0, \mathbf{u}_1).$$

505 We first prove that, given two scalar  $\varepsilon$ -close inputs to a MambaBlock, their deviations do not grow  
 506 exponentially as the number of recurrences increases (Lemma 1). The main result in the paper is  
 507 subsequently proved.

508 **Lemma 1.** *For input  $(\mathbf{x}_0, \mathbf{u}_1)$  to a MambaBlock, small changes  $(\mathbf{x}_0 + \varepsilon, \mathbf{u}_1 + \varepsilon)$  produce deviations  
 509 which are exponentially non-increasing over discrete-time. That is,  $\max |F_\theta^N(\mathbf{x}_0, \mathbf{u}_1) - F_\theta^N(\mathbf{x}_0 + \varepsilon,$   
 510  $\mathbf{u}_1 + \varepsilon)| \in \mathcal{O}(\varepsilon \exp(N\zeta))$ , for some scalar  $\zeta \leq 0$ .*

511 *Proof.* Firstly, we note that within the MambaBlock,  $\mathbf{A}$  is stored in log-space followed by a negative  
 512 exponentiation prior to use. Thus,  $\mathbf{A} \in \mathbb{R}^{d \times d}$ , where is the set of non-positive real numbers.

513 Recall that for the maximum deviation, we have:

$$\max |F_\theta^N(\mathbf{x}_0, \mathbf{u}_1) - F_\theta^N(\mathbf{x}_0 + \varepsilon, \mathbf{u}_1 + \varepsilon)| \in \mathcal{O}(\varepsilon \exp(N\lambda_{\max})).$$

514 where the maximal Lyapunov exponent  $\lambda_{\max}$  is defined as:

$$\lambda_{\max} := \lim_{T \rightarrow \infty} \frac{1}{T} \log \left\| \prod_{t=0}^{T-1} \frac{\partial \mathbf{x}_t}{\partial \mathbf{x}_{t-1}} \right\|_2,$$

515 and  $\|\cdot\|_2$  denotes the spectral norm for matrices.

516 Thus, to complete the proof, it suffices to show that  $\lambda_{\max} \leq 0$ . Recall that  $\mathbf{A}$  and  $\bar{\mathbf{\Delta}}_t$  are diagonal.  
 517 From Equation 5, we thus have

$$\begin{aligned} \lambda_{\max} &= \lim_{T \rightarrow \infty} \frac{1}{T} \log \left\| \prod_{t=0}^{T-1} \frac{\partial \mathbf{x}_t}{\partial \mathbf{x}_{t-1}} \right\|_2 \\ &= \lim_{T \rightarrow \infty} \frac{1}{T} \log \left\| \prod_{t=0}^{T-1} \exp(\bar{\mathbf{\Delta}}_t \mathbf{A}) \right\|_2 \\ &= \lim_{T \rightarrow \infty} \frac{1}{T} \log \left\| \exp \sum_{t=0}^{T-1} (\bar{\mathbf{\Delta}}_t \mathbf{A}) \right\|_2 \end{aligned}$$

518 Let  $i$  be the dimension which corresponds to the output of the spectral norm, i.e.,  $i =$   
 519  $\text{argmax}_{j=1, \dots, d} \{\exp \sum_{t=0}^{T-1} (\bar{\mathbf{\Delta}}_t[j, j] \mathbf{A}[j, j])\}$ . We thus have

$$\begin{aligned} \lambda_{\max} &= \lim_{T \rightarrow \infty} \frac{1}{T} \log \left\| \exp \sum_{t=0}^{T-1} (\bar{\mathbf{\Delta}}_t \mathbf{A}) \right\|_2 \\ &= \lim_{T \rightarrow \infty} \frac{1}{T} \log \exp \sum_{t=0}^{T-1} (\bar{\mathbf{\Delta}}_t[i, i] \mathbf{A}[i, i]) \\ &= \mathbf{A}[i, i] \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=0}^{T-1} \bar{\mathbf{\Delta}}_t[i, i] \end{aligned}$$

520  $\mathbf{A}[i, i]$  is non-positive and  $\lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=0}^T \bar{\Delta}_t[i, i] \geq 0$ , since  $\bar{\Delta}_t[i, i] \in \forall t$ . Thus,  $\lambda_{\max} \leq 0$ .  $\square$

521 **Theorem 2.** Let  $(\mathbf{x}_{t-1}, \mathbf{u}_t)$  be the latent state and input at an arbitrary time  $t \in [1, T]$  within a  
 522 *MambaBlock*. Then small changes  $(\mathbf{x}_{t-1} + \varepsilon, \mathbf{u}_t + \varepsilon)$  produce deviations which are exponentially  
 523 decreasing over discrete-time, i.e.,  $\max |F_{\theta}^N(\mathbf{x}_0, \mathbf{u}_1) - F_{\theta}^N(\mathbf{x}_0 + \varepsilon, \mathbf{u}_1 + \varepsilon)| \in \mathcal{O}(\varepsilon \exp(N\zeta))$ , for  
 524 some scalar  $\zeta \leq 0$ .

525 *Proof.* Let  $\tau(t)$  be a function that maps time values such that  $\tau(t) \in [1, T - t]$  and  $\tau(t) = 1, \tau(t +$   
 526  $1) = 2, \dots, \tau(t + T) = T - t$ . Then  $\mathbf{B}_{\tau(t)}, \mathbf{C}_{\tau(t)}, \Delta_{\tau(t)}$  define a new *MambaBlock* with inputs  
 527  $\mathbf{u}_{\tau(t)}, \dots, \mathbf{u}_{\tau(t+T)}$  and subsequent recurrent states  $\mathbf{x}_{\tau(t)}, \dots, \mathbf{x}_{\tau(t+T)}$ . Applying Lemma 1 to this  
 528 *MambaBlock* with  $(\mathbf{x}_{\tau(t)-1}, \mathbf{u}_{\tau(t)})$  completes the proof.  $\square$

## 529 C Experimental Details

530 All model checkpoints were evaluated on all benchmarks and few-shot settings using the LM  
 531 evaluation harness from Eleuther AI [16], version 0.4.2. Pythia and Mamba Huggingface check-  
 532 points were used for all inference and fine-tuning experiments, e.g., EleutherAI/pythia-160m  
 533 and state-spaces/mamba-130m-hf for the smallest respective models. All fine-tuning experi-  
 534 ments were run using package versions Transformers 4.40.0.dev0, Accelerate 0.28.0, TRL  
 535 0.8.1, PyTorch 2.2.1+cu121, and PEFT 0.10.0.

536 For MPFT, Flash Attention 2.0 [10] via flash\_attn 2.5.7 was used for Pythia mod-  
 537 els. For FP16 and BF16 inference results, Flash Attention 2.0 was used for both Pythia  
 538 and OLMo models. For OLMo results, the 336B-token checkpoint was used by specifying  
 539 revision=step80000-tokens336B.

540 Outside of the OpenHermes hyperparameter search, all Alpaca and OpenHermes fine-tuning exper-  
 541 iments used the following training recipe (adapted from [53]): AdamW\_torch optimizer, cosine  
 542 annealing schedule, no gradient accumulation, maximum norm of 1.0 for gradient clipping, and no  
 543 warmup steps. Training epochs used for all Alpaca and OpenHermes experiments were three and  
 544 one, respectively. For both Pythia and Mamba models, the learning rate and LoRA dimension  $r$  were  
 545 scaled to improve performance of smaller models (per-model values listed in Table 3).

546 For SLL LoRA, targeted Mamba layers were {x\_proj, embeddings, in\_proj, out\_proj};  
 547 x\_proj is the large MambaBlock memory buffer which, when targeted  
 548 by LoRA, regularizes the majority of SSM parameters during fine-tuning  
 549 through weight tying (Theorem 1). Pythia targeted SLL LoRA layers were  
 550 {dense, embed\_in, query\_key\_value, dense\_h\_to\_4h, dense\_4h\_to\_h}, chosen to  
 551 balance performance across model sizes.

552 All experiments in Tables 1 and 2, Figures 1 and 6 were run using a single-GPU Nvidia A10G (24  
 553 GB total memory). For Pythia and Mamba ALL LoRA experiments in Figure 3, all experiments were  
 554 run on an A10G, except for Mamba 2.8B, which exceeded A10G memory capacity and was run on  
 an Nvidia H100 (80 GB total memory).

Table 3: Learning rate and LoRA dimension  $r$  values

Mamba size	Pythia size	learning rate	LoRA $r$
130M	160M	1.0e-5	8
370M	410M	5.0e-5	16
790M	1B	1.0e-6	32
1.4B	1.4B	5.0e-6	64
2.8B	2.8B	5.0e-7	128

555

556 For the hyperparameter search results in Figure 6.1, all experiments were run using 8 H100 GPUs.  
 557 SLL LoRA was used for Mamba and Pythia models. The range of hyperparameter values was as  
 558 follows:

- 559
  - learning rate  $\in \{1e-7, 2e-7, 5e-7, 1e-6, 2e-6, 5e-6, 1e-5, 2e-5, 5e-5, 1e-4, 2e-4, 5e-4, 1e-3, 2e-3, 5e-3\}$
- 560

- 561 • LoRA dimension  $r \in \{16, 32, 64, 128, 256\}$
- 562 • warmup steps  $\in \{0, 1000, 2000\}$

563 All other hyperparameters followed previous experiments.

564 The Alpaca dataset is freely available for download at <https://huggingface.co/datasets/tatsu-lab/alpaca> under open-source license CC-by-NC 4.0. The OpenHermes dataset is freely  
565 available for download at <https://huggingface.co/datasets/teknium/OpenHermes-2.5> un-  
566 der open-source license MIT, Apache 2.0, CC.  
567

## 568 **NeurIPS Paper Checklist**

### 569 **1. Claims**

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

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

573 Justification: All claims made in the abstract and introduction are directly derived from  
574 theoretical and experimental results presented in the main paper.

575 Guidelines:

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

### 585 **2. Limitations**

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

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

588 Justification: Limitations of experimental results are described in the limitations section,  
589 under Discussion.

590 Guidelines:

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

### 617 **3. Theory Assumptions and Proofs**

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

620  
621  
622  
623  
624  
625  
626  
627  
628  
629  
630  
631  
632  
633  
634  
635  
636  
637  
638  
639  
640  
641  
642  
643  
644  
645  
646  
647  
648  
649  
650  
651  
652  
653  
654  
655  
656  
657  
658  
659  
660  
661  
662  
663  
664  
665  
666  
667  
668  
669  
670  
671  
672  
673

Answer: [Yes]

Justification: All theoretical results list any underlying assumptions in the main text and full proofs are available in the supplementary.

Guidelines:

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

#### 4. Experimental Result Reproducibility

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

Answer: [Yes]

Justification: Relevant experimental results are detailed in the main text, with extensive details for all experiments further elaborated upon in the supplementary.

Guidelines:

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

#### 5. Open access to data and code

674 Question: Does the paper provide open access to the data and code, with sufficient instruc-  
675 tions to faithfully reproduce the main experimental results, as described in supplemental  
676 material?

677 Answer: [No]

678 Justification: While we currently answer no, and provide enough detail to reproduce our  
679 experiments, we are actively working towards packaging our code for release. All datasets  
680 are already open source, with licenses listed in the supplementary material.

681 Guidelines:

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

## 701 6. Experimental Setting/Details

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

705 Answer: [Yes]

706 Justification: All datasets are open source, and all experimental hyperparameters are specified  
707 in the paper. All results are fully reproducible with these details.

708 Guidelines:

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

## 714 7. Experiment Statistical Significance

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

717 Answer: [No]

718 Justification: The paper does not report statistical significance.

719 Guidelines:

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

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

## 740 8. Experiments Compute Resources

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

744 Answer: [Yes]

745 Justification: The paper details (at length) the hardware requirements necessary to run each  
746 experiment. Environmental requirements are available as experimental details both in the  
747 main text and supplementary.

748 Guidelines:

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

## 757 9. Code Of Ethics

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

760 Answer: [Yes]

761 Justification: The work detailed in the paper conforms to all aspect of the NeurIPS Code of  
762 Ethics.

763 Guidelines:

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

## 769 10. Broader Impacts

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

772 Answer: [Yes]

773 Justification: The societal impact of this work is addressed in the Discussion section.

774 Guidelines:

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

## 797 11. Safeguards

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

801 Answer: [No]

802 Justification: The work does not aim to release pretrained models or datasets.

803 Guidelines:

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

## 814 12. Licenses for existing assets

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

818 Answer: [Yes]

819 Justification: Extensive lengths were made to cite all original authors for any and all utilized  
820 code/data/work.

821 Guidelines:

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

- 827
- 828
- 829
- 830
- 831
- 832
- 833
- 834
- 835
- 836
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
  - If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets) has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
  - For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
  - If this information is not available online, the authors are encouraged to reach out to the asset's creators.

### 837 13. New Assets

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

840 Answer: [NA]

841 Justification: The paper currently does not release source code. However, as previously  
842 mentioned, we are actively working to remedy this.

843 Guidelines:

- 844
- 845
- 846
- 847
- 848
- 849
- 850
- 851
- The answer NA means that the paper does not release new assets.
  - Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
  - The paper should discuss whether and how consent was obtained from people whose asset is used.
  - At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

### 852 14. Crowdsourcing and Research with Human Subjects

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

856 Answer: [NA]

857 Justification: [TODO]

858 Guidelines:

- 859
- 860
- 861
- 862
- 863
- 864
- 865
- 866
- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
  - Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
  - According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

### 867 15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human 868 Subjects

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

873 Answer: [NA]

874 Justification: [TODO]

875 Guidelines:

- 876
- 877
- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.

878  
879  
880  
881  
882  
883  
884  
885

- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.