

Layer-Wise Analysis in Exploring the Normalization Strategies in Mamba

Peilin Feng^{*1} Yuanshuai Wang^{*1} Yunhao Ni¹ Junle Wang¹ Kui Zhang¹ Wenjun Wu^{1 2 3} Lei Huang^{1 2 3}

Abstract

The Mamba architecture achieves linear time and memory complexity in long-sequence modeling and vision tasks through a dynamic, input-conditioned state transition mechanism and hardware-efficient scan operations. However, as network depth increases, the state space model (SSM) component tends to amplify activation magnitudes during the forward pass, often leading to gradient explosion. To address this, we analyze training stability by tracking (i) the spectral norm of the output projection weights and (ii) the largest eigenvalue of the joint input-output covariance matrix, demonstrating the effectiveness of post-SSM in suppressing activation and gradient scale inflation. From an optimization efficiency perspective, we use K-FAC to approximate the Fisher Information Matrix and show that pre-SSM significantly reduces the condition number of per-layer gradients, thereby accelerating convergence. Furthermore, we propose a composite normalization strategy (BN→SSM→LN), combining BatchNorm at the input and LayerNorm at the output of SSM. We evaluate this strategy across a broad range of benchmarks. Experimental results demonstrate that the composite scheme consistently outperforms single or no normalization in both convergence speed and final accuracy. We hope this work provides both theoretical insights and empirical guidance for normalization in designing SSM-based models.

1. Introduction

Mamba (Gu & Dao, 2023) has attracted significant attention across a broad range of applications, including long-sequence modeling tasks such as speech and audio pro-

^{*}Equal contribution ¹SKLCCSE, Institute of Artificial Intelligence, Beihang University ²Beijing Advanced Innovation Center for Future Blockchain and Privacy Computing, Beihang University ³Hangzhou International Innovation Institute, Beihang University, Hangzhou, China. Correspondence to: Lei Huang <huanlei@buaa.edu.cn>.

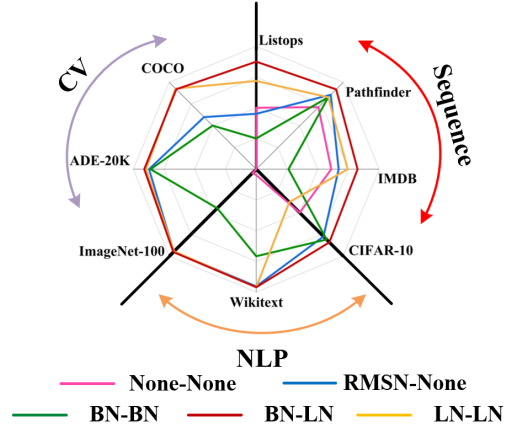


Figure 1. Performance Comparison of Mamba with Composite vs. Single/No Normalization.

cessing (Ren et al., 2025), as well as domains like natural language processing (NLP) and computer vision (CV) (Liu et al., 2024), due to its strong capabilities in capturing long-range dependencies and computational efficiency. However, Mamba faces notable training challenges, where instability can frequently result in training divergence, particularly as model parameters scale up (Dao & Gu, 2024).

The success of deep neural networks depends heavily on improvements in training techniques, particularly the normalization of internal representations (Hinton & Salakhutdinov, 2006; Nair & Hinton, 2010; Kingma & Ba, 2015). As a fundamental component of modern architectures, normalization is widely regarded to stabilize and accelerate training, and further enhance generalization through a scale-invariant analysis, condition analysis, and stochasticity analysis, respectively (Huang, 2022). In particular, prior studies (Zhang & Sennrich, 2019) have demonstrated that Layer Normalization (LN) improves training stability by rescaling hidden activations (Ba, 2016; Zhang & Sennrich, 2019), whereas input Batch Normalization (BN) accelerates convergence by whitening inputs, thereby improving the conditioning of the optimization landscape (Ioffe & Szegedy, 2015).

Despite some recent efforts (Gu & Dao, 2023; Ma et al., 2024; Liu et al., 2024) introducing different normalization

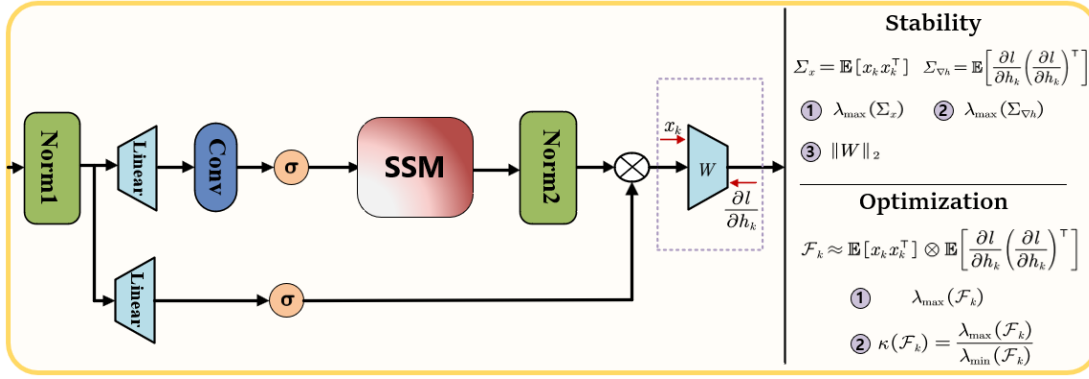


Figure 2. Mamba Normalization Framework. We conduct a layer-wise analysis of training stability and optimization in the Mamba architecture by examining the input and output-gradient of the *out.proj* layer.

layers into the Mamba architecture, such as Root Mean Square Layer Normalization (RMSN) (Zhang & Sennrich, 2019), Layer Normalization (LN) (Ba, 2016) and Group Normalization (GN) (Wu & He, 2018) to improve model performance, there is an absence of systematic analysis on the role of normalization in Mamba and the rationales behind these design choices. To address these challenges, we use the Mamba framework for normalization analysis and choose two critical insertion points for normalization layers, as shown in Figure 2. Norm1, positioned before the input projection layer, which contributes to producing better-conditioned optimization landscapes. And Norm2, positioned after the SSM layer, which ensures scale invariance property to stabilize the training process.

Specifically, we first assess training stability by tracking two key indicators: (i) the spectral norm of the output projection weights and (ii) the maximum singular value of both the layer input covariance matrix and the layer output-gradient covariance matrix. These metrics reflect the stabilizing role of Norm2, which helps enforce scale invariance (Ba, 2016) across layers. After addressing training stability, we shift our focus to optimization efficiency. From this perspective, we analyze the maximum singular value and condition number of the Kronecker-Factored Approximate Curvature (K-FAC) matrix (Huang et al., 2020), which approximates the Fisher Information Matrix (FIM). This analysis reveals how Norm1 substantially enhances optimization conditioning and accelerates convergence behaviors.

Building upon these insights, we propose a composite normalization strategy BN→SSM→LN, employing BN at Norm1 to improve optimization and LN at Norm2 to ensure stability. We evaluate the effectiveness of this design across a diverse set of benchmark tasks, including image classification, object detection, semantic segmentation, long-sequence modeling, and natural language processing. As shown in Figure 1, our BN→SSM→LN configuration con-

sistently outperforms baselines that use either only one normalization method or none at all.

Our contributions are summarized as follows:

- We demonstrate that Norm2 plays a critical role in stabilizing training by suppressing activation and gradient scale explosion, as evidenced by tracking spectral norms and singular values of covariance matrices.
- We show that Norm1 significantly improves optimization efficiency by reducing the condition number of the approximated Fisher Information Matrix K-FAC, thereby accelerating convergence.
- We propose a composite normalization strategy based on these insights, and evaluate its effectiveness across diverse tasks, achieving both faster convergence and higher accuracy compared to using only a single normalization or none.

2. Related Work

In this section, we begin by reviewing the linear State Space Model (SSM) architecture underlying Mamba and its variants. Next, we summarize the role of normalization techniques within deep neural networks. Finally, we introduce key analytical tools—namely the spectral norm of weight matrices, the eigenvalues of joint covariance matrices, and approximations to the Fisher Information Matrix—to assess training stability and optimization efficiency, thereby laying the theoretical groundwork for our analysis.

2.1. Linear State Space Models

State space models (SSMs) achieve linear time complexity and parallelism by either diagonalizing the state transition matrix S4 (Gu et al., 2021) or imposing structural constraints DSS in language modeling and sequence tasks.

Such methods enable the efficient modeling of dependencies spanning tens or even hundreds of thousands of tokens. Building upon this line of research, Gu et al. proposed the Mamba architecture (Gu & Dao, 2023), which introduces a selective mechanism that dynamically adjusts the state transition matrix based on the input, thereby enabling content-aware information propagation. This innovation has inspired several follow-up variants—such as DiMSUM (Phung et al., 2024), Quamba (Chiang et al., 2024), bi-CrossMamba (Wu et al., 2024), Mamba-PTQ (Pierro & Abreu, 2024), CMAMBA (Zeng et al., 2024), and FST-Mamba (Wei et al., 2024) that explore extensions in quantized inference, lightweight design, and hierarchical modeling. Despite these advancements, training Mamba remains challenging, particularly for large-scale models. As shown by (Dao & Gu, 2024), instability during training can lead to divergence, underscoring the need for systematic solutions to stabilize learning dynamics in SSM-based architectures.

2.2. The Role of Normalization in DNNs

The success of deep neural networks (DNNs) heavily relies on training techniques that regulate the distribution of activations (Kingma & Ba, 2015; Ioffe & Szegedy, 2015). Among these, normalization plays a pivotal role. Its benefits primarily lie in enhancing training stability, optimization efficiency, and generalization ability (Huang et al., 2023). For example, Batch Normalization (BN) normalizes activations across both the batch and feature dimensions, mitigating internal covariate shift and improving optimization (Ioffe & Szegedy, 2015; Wang et al., 2022). In contrast, Layer Normalization (LN) performs normalization along the feature dimension within each individual sample, stabilizing hidden state dynamics by fixing the mean and variance of the summed inputs within each layer for individual training examples (Ba, 2016). These have been widely adopted in recurrent neural networks (RNNs) and Transformer-based architectures (Xiong et al., 2020; Shleifer et al., 2022; Han et al., 2021); however, theoretical analysis and investigation remain lacking in the context of Mamba.

2.3. Analysis of Training Stability

The scale-invariant property refers to the characteristic that a network’s behavior or training dynamics remain unaffected when the input or intermediate activations are rescaled (Neyshabur et al., 2015; Wan et al., 2020). This property ensures that gradient descent-based optimization consistently increases the norm of the weights, while the gradients with respect to the weights decrease as the weight norms grow (Huang et al., 2020). Consequently, it stabilizes training by preventing divergence to infinity (Arora et al., 2018). Typical measures of scale invariance include the eigenvalues of the covariance matrix of layer input activations and layer output-gradient covariance matrix, which

together estimate the numerical stability of activations during the training process. Additionally, the spectral norm of the weight matrices is used to assess the scaling behavior of the network itself. Effectively controlling the statistical properties of hidden activations across layers, along with maintaining the spectral norms of weight matrices within reasonable bounds, is a key strategy for achieving scale invariance and ensuring training stability.

2.4. Analysis of Optimization efficiency

After ensuring training stability, another critical consideration is the optimization behavior of deep networks, particularly with respect to convergence speed and final performance. A theoretical foundation for how normalization benefits optimization lies in the approximation of the Fisher Information Matrix (FIM), which characterizes the curvature of the loss landscape. However, computing the full FIM in deep neural networks is computationally prohibitive due to the large number of parameters (Ma et al., 2020). A successful solution to this problem is the Kronecker-Factored Approximate Curvature (K-FAC) method (Martens & Grosse, 2015), which approximates the FIM as a block-diagonal matrix composed of layer-wise sub-FIMs (Ba et al., 2017; Sun & Nielsen, 2017). K-FAC enables a practical and efficient analysis of the optimization landscape at the layer-wise level. Specifically, the maximum eigenvalue and the condition number of each sub-FIM offer insights into gradient magnitudes and ease of optimization, respectively. These quantities effectively reflect the local geometry of the optimization problem across different layers. Building upon this foundation, we adopt a K-FAC-based layer-wise analysis to investigate normalization in the context of the Mamba architecture. This perspective provides a principled and fine-grained approach to understand how input normalization strategies influence the conditioning and optimization behavior during training.

3. Method

To systematically analyze the role of normalization in the Mamba architecture, we adopt the Mamba normalization framework and focus on two critical insertion points: **Norm1**, placed before the *in_proj* layer, and **Norm2**, placed after the structured state space module (SSM). Given that the SSM tends to amplify activation magnitudes during the forward pass, which can accumulate across deep layers and lead to gradient explosion, as shown in Figure 3. We introduce **Norm2** to regulate activation scales thus stabilizing training. To evaluate *stability*, we analyze the spectral norm of the *out_proj* layer weights and the eigenvalues of the layer input and output-gradient covariance matrix. For the *in_proj* side, **Norm1** operates to improve optimization conditions, thereby enhancing learning efficiency, and im-

prove final performance. Accordingly, we compute the maximum eigenvalues and condition number of approximate Fisher Information Matrix K-FAC to assess the influence of **Norm1** on model optimization. Building on the above insights, we propose a composite normalization strategy (BN→SSM→LN) and evaluate its effectiveness across a variety of tasks

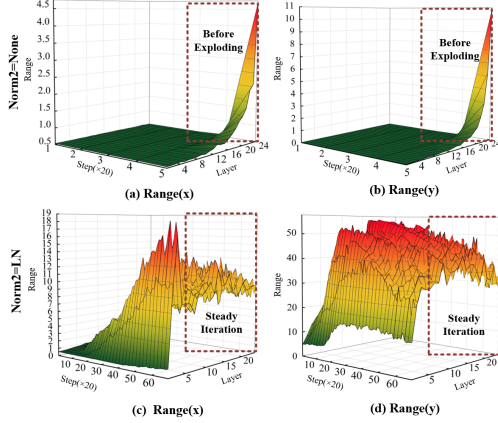


Figure 3. Input and Output Ranges of the SSM Module on the WikiText-103 Dataset (Without vs. With Normalization).

3.1. Preliminaries of Mamba

We first briefly review the core components of the Mamba architecture. As illustrated in Figure 2, let N_1 denote the first normalization layer applied to input x . The main branch then proceeds through a sequence of transformations:

$$f = N_2(\text{SSM}(F_1(N_1(x)))) \quad (1)$$

Here, F_1 represents the main forward path, including a linear projection, depthwise separable convolution, and a SiLU activation. SSM denotes the selective structured state space module, and N_2 is the second normalization layer applied after SSM. Meanwhile, in the parallel branch, the normalized input $N_1(x)$ is processed by a lightweight path F_2 :

$$p = F_2(N_1(x)). \quad (2)$$

F_2 includes a linear projection and a SiLU activation. The outputs of the two branches are combined element-wise and followed by a linear layer:

$$y = L(f \otimes p) \quad (3)$$

Where \otimes denotes element-wise multiplication.

Under this framework, many existing Mamba variants can be seen as specific instances. For example, DiM-SUM (Phung et al., 2024), Quamba (Chiang et al., 2024), bi-CrossMamba (Wu et al., 2024), Mamba-PTQ (Pierro & Abreu, 2024), CMAMBA (Zeng et al., 2024), and FST-Mamba (Wei et al., 2024) adopt $N_1 = \text{RMSN}$; whereas RetinexMamba (Bai et al., 2024b), CU-Mamba (Deng & Gu, 2024), RSMamba (Chen et al., 2024), FMamba (Ma et al., 2024), MLSA4Rec (Su & Huang, 2024), and Zamba (Glorioso et al., 2024) choose $N_1 = \text{LN}$. Unlike the above, DIFFIMP (Gao et al., 2024) and IRSRMamba (Huang et al., 2024) adopt $N_2 = \text{RMSN}$ and $N_2 = \text{LN}$ respectively; Bi-Mamba (Tang et al., 2024), Mamba2 (Bai et al., 2024a), and MambaHSI (Li et al., 2024) explore further variants including $N_2 = \text{Instance Normalization}$ (Ulyanov et al., 2016) and GN.

However, these normalization strategies are mostly designed empirically for specific tasks, and their underlying architectural roles remain unclear. In this study, we explore normalization from two complementary perspectives: training stability and optimization behavior.

3.2. Stability Metrics

To evaluate training stability, we investigate three metrics, the magnitude of layer input (indicated by $\lambda_{\max}(\Sigma_x)$), the magnitude of layer output-gradient (indicated by $\lambda_{\max}(\Sigma_{\nabla_h})$) and the magnitude of *out.proj* weight (indicated by $\|W\|_2$) which ensure that weight magnitudes can grow under gradient descent while gradient norms shrink proportionally—thereby avoiding divergence. To quantify this, we adopt two spectral metrics at the output layer, as shown in Figure 2.

- The spectral norm of the output projection weights, reflecting the scale of activations during training.
- The maximum eigenvalues of the input activation covariance and output gradient covariance matrices, indicating sensitivity to scale perturbations in forward and backward propagation.

Since SSM amplifies activations in forward propagation and accumulates over depth, **Norm2** is used to regulate the activation scale. Prior work has also shown LayerNorm to stabilize training. We therefore conduct controlled experiments comparing $\text{None} \rightarrow \text{SSM} \rightarrow \text{None}$ and $\text{None} \rightarrow \text{SSM} \rightarrow \text{LN}$, and analyze the impact of LN on Mamba training stability using the above metrics.

These metrics jointly indicate whether normalization reduces distortion in activations and gradients across layers, thereby stabilizing overall training dynamics. The subsequent experimental results on Mamba further validate our analysis.

3.3. Optimization Metrics

To evaluate the impact of **Norm1** on model trainability, we analyze the spectral structure of the Fisher Information Matrix (FIM), which characterizes the curvature of the loss landscape, as shown in Figure 2. However, due to memory and compute constraints, directly analyzing the full curvature matrix is infeasible. We instead approximate it using Kronecker-Factored Approximate Curvature (K-FAC) (Huang et al., 2020; Martens & Grosse, 2015). The FIM can be approximated as a block-diagonal matrix:

$$\mathcal{F} \approx \begin{bmatrix} \mathcal{F}_1 & 0 & \cdots & 0 \\ 0 & \mathcal{F}_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \mathcal{F}_L \end{bmatrix}, \quad (4)$$

The k -th block \mathcal{F}_k (or the k -th layer) is approximated as:

$$\mathcal{F}_k \approx \mathbb{E}[x_k x_k^\top] \otimes \mathbb{E} \left[\frac{\partial l}{\partial h_k} \left(\frac{\partial l}{\partial h_k} \right)^\top \right] \quad (5)$$

Where, x_k is the input to the k -th layer, and $\frac{\partial l}{\partial h_k}$ is the output gradient.

Since **Norm1** is intended to improve optimization conditions and convergence, we conduct experiments comparing $\text{None} \rightarrow \text{SSM} \rightarrow \text{LN}$ and $\text{BN} \rightarrow \text{SSM} \rightarrow \text{LN}$. We analyze the impact of BN using the condition number $\kappa(\mathcal{F}_k)$ ¹:

$$\kappa(\mathcal{F}_k) = \frac{\lambda_{\max}(\mathcal{F}_k)}{\lambda_{\min}(\mathcal{F}_k)} \quad (6)$$

A lower $\kappa(\mathcal{F}_k)$ implies better conditioning and more efficient gradient-based optimization. A higher condition number indicates ill-conditioning and potential convergence challenges. The experimental results on the Mamba architecture presented later also support this analysis.

3.4. Composite Strategy

In DNNs, Batch Normalization (BN) facilitates faster training (Ioffe & Szegedy, 2015), while Layer Normalization (LN) helps stabilize the training process (Ba, 2016). This motivates our proposed composite normalization strategy: **BN**→**SSM**→**LN**, as shown in Figures 4, which aims to achieve both goals simultaneously. The combined normalization can be expressed as:

$$f = \text{LN}(\text{SSM}(F_1(\text{BN}(x)))). \quad (7)$$

¹The general condition number with respect to the percentage is defined as: $\kappa_p = \frac{\lambda_{\max}}{\lambda_p}$ where λ_p is the p -th eigenvalue (in descending order). This measure provides a better characterization of over-parameterized models.

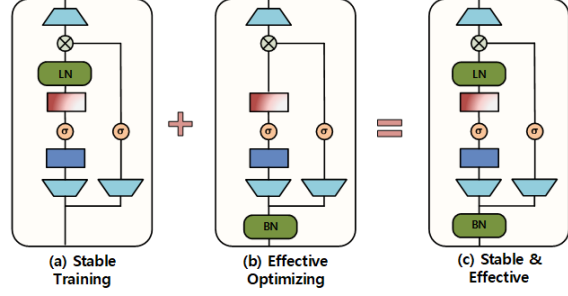


Figure 4. Composite Normalization Strategy: Integrating Stability and Optimization Configurations

Where a Batch Normalization (BN) is first applied to the input x , followed by F_1 , after which the result is passed through the state space model (SSM) and finally normalized by a Layer Normalization (LN). To evaluate the effectiveness of the proposed combined normalization, we compare it against baseline normalization methods, which are $\text{None} \rightarrow \text{SSM} \rightarrow \text{None}$, $\text{RMSN} \rightarrow \text{SSM} \rightarrow \text{None}$, $\text{LN} \rightarrow \text{SSM} \rightarrow \text{LN}$, and $\text{BN} \rightarrow \text{SSM} \rightarrow \text{BN}$.

In the following section, we conduct experiments across diverse tasks to validate the effectiveness and generalizability of the proposed method.

4. Experiments

In this section, we first introduce the datasets and experimental settings used to evaluate the impact of normalization on the Mamba architecture across vision, natural language processing, and sequential tasks. Next, we analyze the normalization results in language modeling and image classification tasks using output-layer weight norms, eigenvalues of input-gradient covariance matrices, and K-FAC condition numbers. Finally, we conduct comparison experiments on our proposed composite BN and LN normalization strategy across various tasks to verify its generalizability.

4.1. Experiment Settings

Baselines We select the vanilla Mamba architecture, which adopts RMSNorm-None as the normalization configuration, and the widely used VMamba architecture, which employs LN-LN as its normalization setup. Both serve as baselines for comparison.

Datasets We use a range of datasets to evaluate the performance across different tasks. For stability analysis and optimization analysis, we utilize WikiText-103 (Merity et al., 2016), a widely-used dataset for language modeling, and ImageNet-100 (Ima, 2019), a subset of ImageNet for image classification. For generalization verification, we evaluate the combined normalization strategy across a vari-

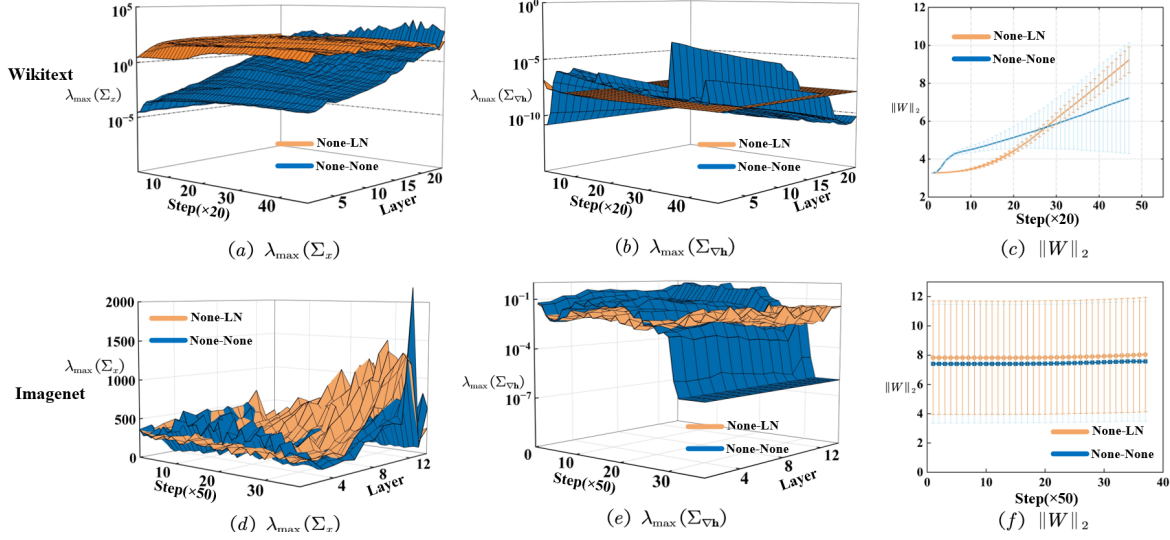


Figure 5. Analysis of layer input magnitude, output gradient magnitude, and weight norm. Yellow indicates None→LN, and blue indicates None→None. Subfigures (a), (b), and (c) illustrate the variations in stability metrics on the WikiText-103 dataset, while (d), (e), and (f) present the corresponding results on the ImageNet-100 dataset.

ety of benchmark datasets, including sequence tasks from the LRA Benchmark (Tay et al., 2021), NLP tasks with WikiText-103, and computer vision tasks such as ImageNet-100, COCO (Lin et al., 2014) and ADE-20K (Zhou et al., 2019). The dataset and experimental configurations are described in detail in the Appendix 1.

4.2. Stability Analysis

We begin by examining the impact of normalization strategies on training stability. Following the setup described in the Method section, we compare two configurations: None→SSM→LN and None→SSM→None, which correspond to applying LayerNorm after the SSM versus no normalization at all.

On the WikiText-103 and ImageNet-100 dataset, we track the spectral norm of output projection weights, as well as the maximum eigenvalues of the input activation covariance and output gradient covariance matrices across Mamba layers, as shown in Figure 5. The results are summarized below:

- **Weight Norms:** Under the None→SSM→None configuration, the weight norms of deeper layers (e.g., layer 20) increase significantly, far exceeding earlier layers. This results in gradient explosion and even training divergence. In contrast, with LayerNorm (None→SSM→LN), the norm trends remain consistent across layers, and gradients maintain scale invariance, enabling smoother training, as shown in Figures 5(a) and (d).
- **Output Gradient Covariance Eigenvalues:** Com-

pared to the None→SSM→None, None→SSM→LN exhibits more consistent gradient eigenvalue distributions and reduced fluctuations during training, suggesting smoother gradient flow, as shown in Figures 5(b) and (e).

- **Input Covariance Eigenvalues:** The None→SSM→LN setup maintains consistent and relatively high eigenvalues across layers, with minimal variation over training iterations, indicating effective suppression of forward-pass scale perturbation. Without normalization, inter-layer eigenvalue differences are large, reducing numerical stability, as shown in Figures 5(c) and (f).

These results confirm that output-side normalization (Norm2) significantly suppresses activation and gradient explosion, thereby improving the training stability of deep Mamba networks. This also validates the theoretical insights in Section 3, where Norm2 was shown to alleviate scale inflation caused by the SSM.

4.3. Optimization Analysis

We further investigate the effect of input-side normalization (Norm1) on optimization efficiency by applying BN before the SSM. We compare BN→SSM→LN against None→SSM→LN, using the maximum eigenvalue and condition number of the K-FAC-approximated Fisher Information Matrix as evaluation metrics.

Results on the WikiText-103 and ImageNet-100 datasets show that:

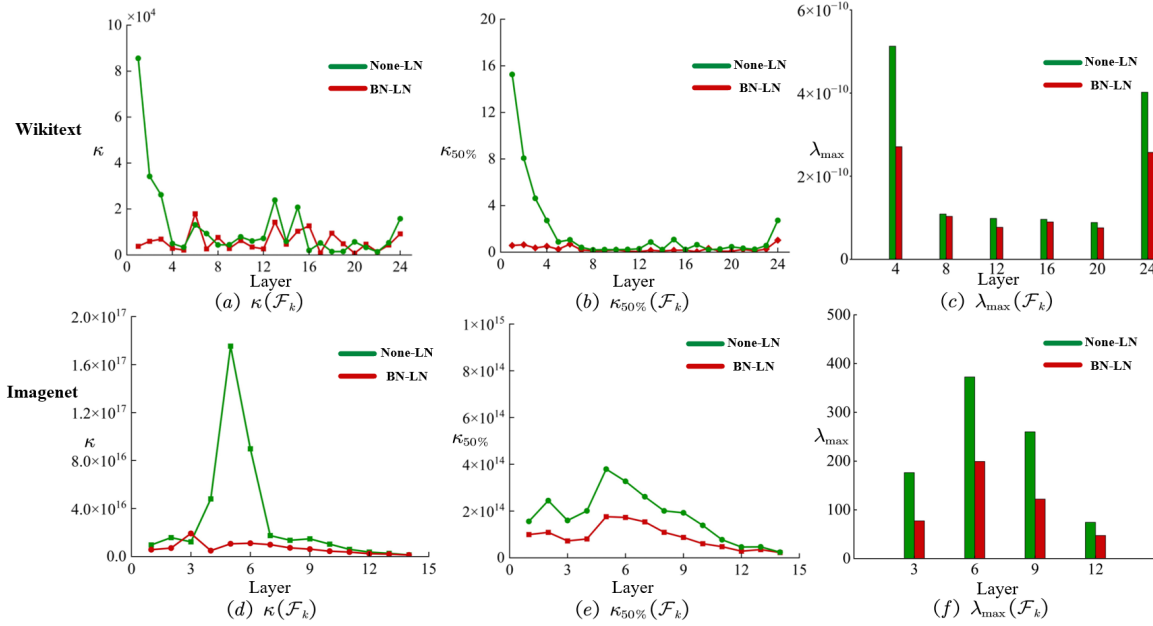


Figure 6. Analysis of the condition of K-FAC (indicated by κ and $\kappa_{50\%}$) and magnitude of K-FAC (indicated by $\lambda_{\max}(\mathcal{F}_k)$). The green line represents None \rightarrow LN, and red line indicates BN \rightarrow LN. Subfigures (a), (b), and (c) illustrate the corresponding optimization metric changes on the WikiText-103 dataset, while (d), (e), and (f) are on the ImageNet-100 dataset.

- **K-FAC Condition Number:** Across 100% and 50% thresholds, the condition numbers under BN \rightarrow SSM \rightarrow LN (with BN) are significantly lower than those without BN (None \rightarrow SSM \rightarrow LN), indicating faster gradient convergence and improved training efficiency, as shown in Figures 6(a) and (d).
- **Convergence Performance:** Compared to the None \rightarrow SSM \rightarrow LN (without BN), BN \rightarrow SSM \rightarrow LN (with BN) helps the Mamba reach lower training loss and better generalization performance more rapidly, as shown in Figures 6(b) and (e).
- **K-FAC Maximum Eigenvalue:** The K-FAC Maximum Eigenvalue under BN \rightarrow SSM \rightarrow LN are lower than under None \rightarrow SSM \rightarrow LN, suggesting better alignment in parameter update directions and a smoother optimization landscape, as shown in Figures 6(c) and (f).

These results indicate that input-side BN not only accelerates convergence but also improves numerical conditioning during optimization, thereby enhancing the trainability of Mamba models.

4.4. Validation of BN-LN Composite Normalization

Building on the above theoretical and empirical analyses, we propose the composite normalization strategy BN \rightarrow SSM \rightarrow LN, and conduct systematic comparisons across tasks including vision classification, segmentation, and

reasoning, sequence modeling, natural language processing,. The datasets include ImageNet-100, COCO, ADE-20K, Pathfinder, ListOps, CIFAR-10, IMDB (Text), and WikiText-103. Results are summarized in Tables 1, 2, and 3, respectively.

Table 1. Results of different normalization strategies on sequence tasks. Configurations that result in divergent (NaN) losses during training are marked with an asterisk (*).

| Method | ListOps | CIFAR | Pathfinder |
|--|--------------|--------------|--------------|
| None \rightarrow None | 38.61* | 56.4* | 49.95 |
| RMSN \rightarrow None | 39.51 | 62.74 | 51.00 |
| BN \rightarrow BN | 37.50* | 63.09 | 50.80* |
| LN \rightarrow LN | 42.18 | 58.80 | 50.80 |
| BN\rightarrowLN (Ours) | 43.75 | 63.41 | 51.43 |

It can be observed that single-use BN or LN strategies lead to unstable or divergent behavior in certain tasks. In contrast, the BN-LN composite strategy not only significantly accelerates convergence but also achieves the best (or even state-of-the-art) performance across all evaluated tasks. Particularly in deeper Mamba models, BN-LN effectively balances optimization speed and training stability, demonstrating stronger generalization.

Moreover, the evaluation metrics curves during training are shown in Figures 7 and 8, respectively. These figures also demonstrate that combined normalization leads to faster

Table 2. Results of different normalization strategies on NLP task WikiText-103. Configurations that result in divergent (NaN) losses during training are marked with an asterisk (*).

| Method | WikiText-103 | IMDB |
|---------------------|--------------|--------------|
| None→None | 201.07* | 77.2* |
| RMSN→None | 28.9 | 78.40 |
| BN→BN | 201.3* | 70.24 |
| LN→LN | 27.59 | 79.87 |
| BN→LN (Ours) | 27.57 | 81.48 |

Table 3. Results of different normalization strategies on visual tasks. Configurations that result in divergent (NaN) losses during training are marked with an asterisk (*).

| Method | ImageNet100 | COCO | ADE20K |
|---------------------|--------------|-------------|--------------|
| None→None | 23.08* | 0* | 0* |
| RMSN→None | 86.34 | 24.2* | 26.17 |
| BN→BN | 51.38* | 20.1 | 25.78 |
| LN→LN | 86.16 | 34.5 | 26.92 |
| BN→LN (Ours) | 86.72 | 34.9 | 27.32 |

convergence. For example, in the segmentation task, the combined normalization consistently outperforms single normalization methods in terms of accuracy and reaches the highest accuracy earlier during training.

To further evaluate the effectiveness of different normalization strategies in accelerating convergence, we replaced Norm1 with several commonly used normalization methods, as shown in Figure 9. It can be seen that the combined normalization configuration of BN+LN not only maintains high final accuracy but also achieves the fastest convergence, making it the optimal choice for efficient training.

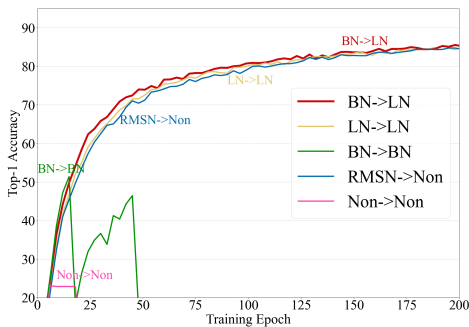


Figure 7. Training stability and convergence curves on ImageNet classification

5. Conclusion

In this paper, we investigate the training stability and optimization convergence of normalization in the Mamba archi-

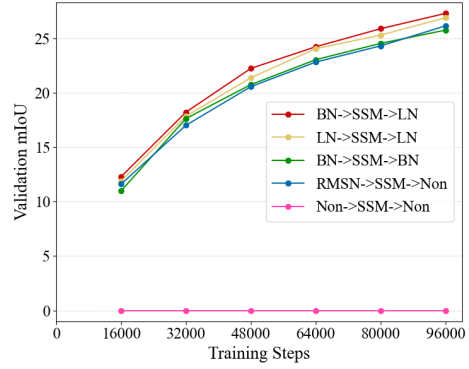


Figure 8. Training stability and convergence on ADE-20K segmentation

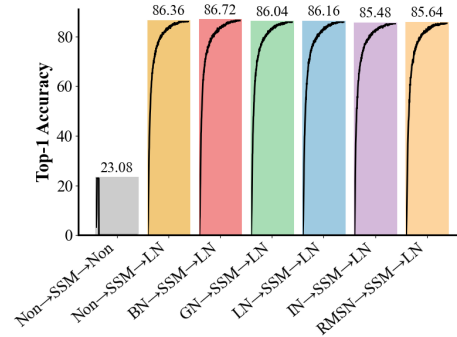


Figure 9. Comparison with various Norm1 methods.

ture. In terms of training stability, by monitoring the spectral norm of the output projection weights and the largest eigenvalue of the joint input-output covariance matrix, we demonstrate that post-SSM LayerNorm is crucial for suppressing activation and gradient magnification and thus preventing gradient explosion in deep networks. With respect to optimization efficiency, through condition-number estimates of a K-FAC-approximated Fisher Information Matrix, we show that pre-SSM BatchNorm substantially improves numerical conditioning, accelerating gradient convergence and training speed. Finally, across tasks such as language modeling, image classification, and semantic segmentation, our composite strategy not only converges more rapidly but also outperforms baselines. However, this study is limited by its focus on only two normalization methods and fixed insertion positions, which may restrict scalability to deeper networks and larger-scale tasks. Future work will explore broader normalization variants, automated placement strategies, and extend the framework to more complex architectures and large-scale settings to enhance generality and performance.

Impact Statement

This paper presents work whose goal is to advance the field of Deep Learning. There are many potential social consequences of our work, none which feel must be specifically highlighted here.

References

- Imagenet-100: A subset of imagenet-1k with 100 randomly selected classes. <https://github.com/HobbitLong/CMC>, 2019. HuggingFace dataset “clane9/imagenet-100”, derived from ImageNet-1K.
- Arora, S., Li, Z., and Lyu, K. Theoretical analysis of auto rate-tuning by batch normalization. *arXiv preprint arXiv:1812.03981*, 2018.
- Ba, J., Grosse, R., and Martens, J. Distributed second-order optimization using kronecker-factored approximations. In *International conference on learning representations*, 2017.
- Ba, J. L. Layer normalization. *ArXiv preprint*, abs/1607.06450, 2016. URL <https://arxiv.org/abs/1607.06450>.
- Bai, J., Fang, Y., Wang, J., and Zhang, X. A two-stage band-split mamba-2 network for music separation. *ArXiv preprint*, abs/2409.06245, 2024a. URL <https://arxiv.org/abs/2409.06245>.
- Bai, J., Yin, Y., He, Q., Li, Y., and Zhang, X. Retinexmamba: Retinex-based mamba for low-light image enhancement. *ArXiv preprint*, abs/2405.03349, 2024b. URL <https://arxiv.org/abs/2405.03349>.
- Chen, K., Chen, B., Liu, C., Li, W., Zou, Z., and Shi, Z. Rsmamba: Remote sensing image classification with state space model. *ArXiv preprint*, abs/2403.19654, 2024. URL <https://arxiv.org/abs/2403.19654>.
- Chiang, H.-Y., Chang, C.-C., Frumkin, N., Wu, K.-C., and Marculescu, D. Quamba: A post-training quantization recipe for selective state space models. *ArXiv preprint*, abs/2410.13229, 2024. URL <https://arxiv.org/abs/2410.13229>.
- Dao, T. and Gu, A. Transformers are ssms: Generalized models and efficient algorithms through structured state space duality. *ICML*, 2024.
- Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., and Fei-Fei, L. Imagenet: A large-scale hierarchical image database. *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pp. 248–255, 2009. doi: 10.1109/CVPR.2009.5206848.
- Deng, R. and Gu, T. Cu-mamba: Selective state space models with channel learning for image restoration. *ArXiv preprint*, abs/2404.11778, 2024. URL <https://arxiv.org/abs/2404.11778>.
- Desjardins, G., Simonyan, K., Pascanu, R., and Kavukcuoglu, K. Natural neural networks. In Cortes, C., Lawrence, N. D., Lee, D. D., Sugiyama, M., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems 28: Annual Conference on Neural Information Processing Systems 2015, December 7-12, 2015, Montreal, Quebec, Canada*, pp. 2071–2079, 2015. URL <https://proceedings.neurips.cc/paper/2015/hash/2de5d16682c3c35007e4e92982f1a2ba-Abstract.html>.
- Fasina, O., Huguet, G., Tong, A., Zhang, Y., Wolf, G., Nickel, M., Adelstein, I., and Krishnaswamy, S. Neural fim for learning fisher information metrics from point cloud data. In *International Conference on Machine Learning*, pp. 9814–9826. PMLR, 2023.
- Gao, H., Shen, W., Qiu, X., Xu, R., Hu, J., and Yang, B. Diffimp: Efficient diffusion model for probabilistic time series imputation with bidirectional mamba backbone. *ArXiv preprint*, abs/2410.13338, 2024. URL <https://arxiv.org/abs/2410.13338>.
- Glorioso, P., Anthony, Q., Tokpanov, Y., Whittington, J., Pilault, J., Ibrahim, A., and Millidge, B. Zamba: A compact 7b ssm hybrid model. *ArXiv preprint*, abs/2405.16712, 2024. URL <https://arxiv.org/abs/2405.16712>.
- Gu, A. and Dao, T. Mamba: Linear-time sequence modeling with selective state spaces. *ArXiv preprint*, abs/2312.00752, 2023. URL <https://arxiv.org/abs/2312.00752>.
- Gu, A., Goel, K., and Ré, C. Efficiently modeling long sequences with structured state spaces. *ArXiv preprint*, abs/2111.00396, 2021. URL <https://arxiv.org/abs/2111.00396>.
- Han, K., Xiao, A., Wu, E., Guo, J., Xu, C., and Wang, Y. Transformer in transformer. *Advances in neural information processing systems*, 34:15908–15919, 2021.
- Hinton, G. E. and Salakhutdinov, R. R. Reducing the dimensionality of data with neural networks. *science*, 313 (5786):504–507, 2006.
- Huang, L. *Normalization Techniques in Deep Learning*. Springer, 2022.

- Huang, L., Qin, J., Liu, L., Zhu, F., and Shao, L. Layer-wise conditioning analysis in exploring the learning dynamics of dnns. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part II 16*, pp. 384–401. Springer, 2020.
- Huang, L., Qin, J., Zhou, Y., Zhu, F., Liu, L., and Shao, L. Normalization techniques in training dnns: Methodology, analysis and application. *IEEE transactions on pattern analysis and machine intelligence*, 45(8):10173–10196, 2023.
- Huang, Y., Miyazaki, T., Liu, X., and Omachi, S. Irsr-mamba: Infrared image super-resolution via mamba-based wavelet transform feature modulation model. *ArXiv preprint*, abs/2405.09873, 2024. URL <https://arxiv.org/abs/2405.09873>.
- Ioffe, S. and Szegedy, C. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In Bach, F. R. and Blei, D. M. (eds.), *Proceedings of the 32nd International Conference on Machine Learning, ICML 2015, Lille, France, 6-11 July 2015*, volume 37 of *JMLR Workshop and Conference Proceedings*, pp. 448–456. JMLR.org, 2015. URL <http://proceedings.mlr.press/v37/loffel15.html>.
- Kim, J., Lee, B., Park, C., Oh, Y., Kim, B., Yoo, T., Shin, S., Han, D., Shin, J., and Yoo, K. M. Peri-In: Revisiting layer normalization in the transformer architecture. *arXiv e-prints*, pp. arXiv–2502, 2025.
- Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization. In Bengio, Y. and LeCun, Y. (eds.), *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015. URL <http://arxiv.org/abs/1412.6980>.
- Krizhevsky, A., Hinton, G., et al. Learning multiple layers of features from tiny images. 1(4):7, 2009.
- Li, Y., Luo, Y., Zhang, L., Wang, Z., and Du, B. Mambahsi: Spatial-spectral mamba for hyperspectral image classification. *IEEE Transactions on Geoscience and Remote Sensing*, 2024.
- Lin, T.-Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., and Zitnick, C. L. Microsoft coco: Common objects in context. *Computer Vision–ECCV 2014*, pp. 740–755, 2014. doi: 10.1007/978-3-319-10602-1_48.
- Linsley, D., Ashok, A., Govindarajan, L., Liu, R., and Serre, T. Learning long-range spatial dependencies with horizontal gated recurrent units. *Advances in Neural Information Processing Systems*, 31, 2018.
- Liu, Y., Tian, Y., Zhao, Y., Yu, H., Xie, L., Wang, Y., Ye, Q., and Liu, Y. Vmamba: Visual state space model 2024. *CVPR*, 2024.
- Loshchilov, I. and Hutter, F. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.
- Ma, L., Montague, G., Ye, J., Yao, Z., Gholami, A., Keutzer, K., and Mahoney, M. Inefficiency of k-fac for large batch size training. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pp. 5053–5060, 2020.
- Ma, S., Kang, Y., Bai, P., and Zhao, Y.-B. Fmamba: Mamba based on fast-attention for multivariate time-series forecasting. *ArXiv preprint*, abs/2407.14814, 2024. URL <https://arxiv.org/abs/2407.14814>.
- Maas, A. L., Daly, R. E., Pham, P. T., Huang, D., Ng, A. Y., and Potts, C. Learning word vectors for sentiment analysis. In *Proceedings of the 49th annual meeting of the association for computational linguistics: Human language technologies*, pp. 142–150. Association for Computational Linguistics, 2011.
- Martens, J. and Grosse, R. Optimizing neural networks with kronecker-factored approximate curvature. In *International conference on machine learning*, pp. 2408–2417. PMLR, 2015.
- Merity, S., Xiong, C., Bradbury, J., and Socher, R. Pointer sentinel mixture models, 2016.
- Nair, V. and Hinton, G. E. Rectified linear units improve restricted boltzmann machines. In *Proceedings of the 27th international conference on machine learning (ICML-10)*, pp. 807–814, 2010.
- Nangia, N. and Bowman, S. R. Listops: A diagnostic dataset for latent tree learning. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Student Research Workshop*, pp. 92–99, 2018.
- Neyshabur, B., Tomioka, R., Salakhutdinov, R., and Srebro, N. Data-dependent path normalization in neural networks. *arXiv preprint arXiv:1511.06747*, 2015.
- Oczkowski, W. J. and Barreca, S. Neural network modeling accurately predicts the functional outcome of stroke survivors with moderate disabilities. *Archives of physical medicine and rehabilitation*, 78(4):340–345, 1997.
- Phung, H., Dao, Q., Dao, T., Phan, H., Metaxas, D., and Tran, A. Dimsum: Diffusion mamba—a scalable and unified spatial-frequency method for image generation. *ArXiv preprint*, abs/2411.04168, 2024. URL <https://arxiv.org/abs/2411.04168>.

- Pierro, A. and Abreu, S. Mamba-ptq: Outlier channels in recurrent large language models. *ArXiv preprint*, abs/2407.12397, 2024. URL <https://arxiv.org/abs/2407.12397>.
- Ren, W., Wu, H., Lin, Y.-C., Chen, X., Chao, R., Hung, K.-H., Li, Y.-J., Ting, W.-Y., Wang, H.-M., and Tsao, Y. Leveraging joint spectral and spatial learning with mamba for multichannel speech enhancement. In *ICASSP 2025-2025 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 1–5. IEEE, 2025.
- Saarinen, S., Bramley, R., and Cybenko, G. Ill-conditioning in neural network training problems. *SIAM Journal on Scientific Computing*, 14(3):693–714, 1993.
- Shleifer, S., Weston, J., and Ott, M. Normformer: Improved transformer pretraining with extra normalization. *ICLR*, 2022.
- Su, J. and Huang, Z. Mls4rec: Mamba combined with low-rank decomposed self-attention for sequential recommendation. *ArXiv preprint*, abs/2407.13135, 2024. URL <https://arxiv.org/abs/2407.13135>.
- Sun, K. and Nielsen, F. Relative fisher information and natural gradient for learning large modular models. In *International Conference on Machine Learning*, pp. 3289–3298. PMLR, 2017.
- Tang, S., Ma, L., Li, H., Sun, M., and Shen, Z. Bi-mamba: Towards accurate 1-bit state space models. *ArXiv preprint*, abs/2411.11843, 2024. URL <https://arxiv.org/abs/2411.11843>.
- Tay, Y., Dehghani, M., Bahri, D., and Metzler, D. Long range arena: A benchmark for efficient transformers. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2021.
- Ulyanov, D., Vedaldi, A., and Lempitsky, V. Instance normalization: The missing ingredient for fast stylization. *arXiv preprint arXiv:1607.08022*, 2016.
- Wan, R., Zhu, Z., Zhang, X., and Sun, J. Spherical motion dynamics: Learning dynamics of neural network with normalization, weight decay, and sgd. *arXiv preprint arXiv:2006.08419*, 2020.
- Wang, J., Wu, J., and Huang, L. Understanding the failure of batch normalization for transformers in nlp. *Advances in Neural Information Processing Systems*, 35:37617–37630, 2022.
- Wei, Y., Abrol, A., Hassanzadeh, R., and Calhoun, V. Hierarchical spatio-temporal state-space modeling for fmri analysis. *ArXiv preprint*, abs/2408.13074, 2024. URL <https://arxiv.org/abs/2408.13074>.
- Wu, D., Wang, Y., Wu, X., and Qu, T. Cross-attention inspired selective state space models for target sound extraction. *ArXiv preprint*, abs/2409.04803, 2024. URL <https://arxiv.org/abs/2409.04803>.
- Wu, Y. and He, K. Group normalization. In *Proceedings of the European conference on computer vision (ECCV)*, pp. 3–19, 2018.
- Xiong, R., Yang, Y., He, D., Zheng, K., Zheng, S., Xing, C., Zhang, H., Lan, Y., Wang, L., and Liu, T. On layer normalization in the transformer architecture. In *International conference on machine learning*, pp. 10524–10533. PMLR, 2020.
- Xu, J., Sun, X., Zhang, Z., Zhao, G., and Lin, J. Understanding and improving layer normalization. *Advances in neural information processing systems*, 32, 2019.
- Zeng, C., Liu, Z., Zheng, G., and Kong, L. C-mamba: Channel correlation enhanced state space models for multivariate time series forecasting. *ArXiv preprint*, abs/2406.05316, 2024. URL <https://arxiv.org/abs/2406.05316>.
- Zhang, B. and Sennrich, R. Root mean square layer normalization. In Wallach, H. M., Larochelle, H., Beygelzimer, A., d’Alché-Buc, F., Fox, E. B., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pp. 12360–12371, 2019. URL <https://proceedings.neurips.cc/paper/2019/hash/1e8a19426224ca89e83cef47f1e7f53b-Abstract.html>.
- Zhou, B., Zhao, H., Puig, X., Fidler, S., Barriuso, A., and Torralba, A. Semantic understanding of scenes through the ade20k dataset. *International Journal of Computer Vision*, 127:302–321, 2019. doi: 10.1007/s11263-018-1140-0.

Appendix

In this appendix, we provide additional details that could not be included in the main paper due to limited space, which comprises the details of our experiment settings, theoretical backgrounds, empirical experiments, and theoretical analysis. We discuss:

- Datasets and experiment settings.
- Supplementary Theoretical Backgrounds
- Empirical evidence for the role of **Norm2** in enhancing **training stability**.
- Empirical evidence for the role of **Norm1** in improving **optimization condition**.
- Discussion of other composite Strategy cases.

A. Datasets and Experiment Settings

We conduct experiments on benchmark datasets spanning three domains: sequential modeling, natural language processing (NLP), and computer vision (CV).

A.1. NLP Task

Dataset WikiText-103 (Merity et al., 2016) is a large-scale English word-level language modeling benchmark consisting of 28,475 high-quality Wikipedia articles. It retains original casing, punctuation, and numerical content, with the training set comprising approximately 103 million words and a vocabulary of over 260,000 tokens. The validation and test sets each contain 60 full articles. Notably, the dataset preserves paragraph continuity within articles, making it well-suited for evaluating a model’s ability to capture long-range dependencies across thousands of tokens.

Experiment Setting Our Mamba-based language model comprises 24 layers with a hidden dimension of 768, totaling approximately 125 million parameters. The model adopts the Mamba1 state space architecture, without employing RMS normalization or tying input and output embeddings. We trained the 24-layer Mamba for 150 epochs using Distributed Data Parallel (DDP) across 8 GPUs, with a global batch size of 128. We use the AdamW optimizer (Loshchilov & Hutter, 2017) with a peak learning rate of 1.5×10^{-3} and a weight decay of 0.25. The learning rate follows a cosine annealing schedule with 1% linear warm-up steps, starting from 1×10^{-6} and decaying to 10% of the peak value. Gradient clipping is applied with a maximum norm of 1.0. All computations are performed using FP32 precision.

A.2. Sequential Modeling Benchmark

In Long range arena(LRA) benchmark (Tay et al., 2021), we use a 8-layer Mamba1-based sequence classification model with a hidden size of 128 and approximately 1.4M parameters. It uses a state dimension of 64, kernel size of 4, expansion factor of 2, and no normalization layers. Positional encodings are added to capture sequence order.

Dataset ListOps (Nangia & Bowman, 2018) contains approximately 90,000 training samples, 10,000 validation samples, and 10,000 test samples, totaling around 110,000 prefix-style arithmetic expressions with nested operations. Each sequence has an average length of 130 tokens, with some exceeding 200 tokens. The task requires outputting a single integer between 0–9, with operators such as MAX, MIN, MED, and SUM_MOD (SM). This dataset evaluates a model’s ability to reason over long-distance dependencies and recursive tree structures, using accuracy as the evaluation metric.

Experiment Setting The model is trained on LISTOPS for 40 epochs using AdamW (learning rate 1×10^{-4} , weight decay 0.05) with a constant schedule and 2,000 warm-up steps. Training is performed with DDP over 8 GPUs, batch size 64, and no gradient clipping. Inputs are padded to 2,048 tokens with a vocabulary size of 18. End-of-sequence tokens are appended, and outputs are aggregated via pooling with length-aware processing.

Dataset IMDB (Maas et al., 2011) is a sentiment analysis dataset consisting of English movie reviews. It provides 25,000 labeled training samples, 25,000 labeled test samples, and an additional 50,000 unlabeled reviews for semi-supervised learning. Review lengths range from 200 to 1,000 words, and labels are binary: positive or negative. The dataset features a balanced sentiment distribution and linguistic diversity, including slang, negation, and sarcasm, making it a standard benchmark for assessing a model’s capacity to capture fine-grained sentiment in long-form text. Accuracy is used as the evaluation metric.

Experiment Setting The model is trained on the IMDB dataset for 65 epochs with a batch size of 32. We use the AdamW optimizer with a learning rate of 1×10^{-4} , a weight decay of 0.1, and a constant learning rate schedule with 2,000 linear warm-up steps (approximately one epoch). Input sequences are tokenized at the character level using a minimum frequency threshold of 15, yielding a vocabulary of 135 characters. Sequences are padded to a maximum length of 4,096 characters, and end-of-sequence tokens are appended. The final outputs are computed via pooling-based sequence classification with length-aware aggregation to effectively handle variable-length movie reviews.

Dataset CIFAR-10 (Krizhevsky et al., 2009) consists of 50,000 training images and 10,000 test images, totaling 60,000 32×32 RGB images across 10 common categories: airplane, automobile, bird, cat, deer, dog, frog, horse, ship, and truck. The low resolution and cluttered backgrounds demand that models learn discriminative features from limited pixels. Each class contains 5,000 training samples. The dataset is widely used to assess image classification capabilities, with accuracy serving as the evaluation metric.

Experiment Setting The model processes CIFAR-10 images converted to grayscale and serialized into 1,024-token sequences (32×32 pixels), without any data augmentation. Training is conducted for 150 epochs with a batch size of 50. We use the AdamW optimizer with a learning rate of 1×10^{-3} , weight decay of 0.1, and β parameters set to (0.9, 0.95). The learning rate follows a cosine annealing schedule with 2,000 linear warm-up steps. Gradient clipping is applied with a maximum norm of 1.0 to further stabilize training. Each pixel is treated as a discrete token, allowing sequence modeling techniques to be applied to vision tasks through this serialization-based approach.

Dataset Pathfinder (Linsley et al., 2018) is a visual reasoning benchmark designed to assess topological reasoning. The task is to determine whether two marked circles in a binary image are connected by a single continuous path. The dataset includes approximately 100,000 training images and 20,000 validation/test images, with each image sized at 64×64 pixels. As the number of path segments increases, so does task difficulty. Since no semantic cues are available, models must rely on global receptive fields and spatial reasoning. Evaluation is based on accuracy.

Experiment Setting The model processes PATHFINDER images converted to grayscale and serialized into sequences of varying lengths: 1,024 tokens (32×32), 4,096 tokens (64×64), and 65,536 tokens (256×256). Training is conducted for 200 epochs with a batch size of 32. We employ the AdamW optimizer with learning rates of 1×10^{-4} for lower resolutions and 1×10^{-3} for the 256×256 setting, using weight decay values between 0.05 and 0.1, and β parameters set to (0.9, 0.95). The learning rate follows either a constant or cosine annealing schedule with 5,000 linear warm-up steps. Gradient clipping with a maximum norm of 1.0 is applied for higher-resolution inputs. Each pixel is treated as a discrete token in the serialized sequence, allowing us to evaluate the model’s capacity to capture long-range dependencies across different input lengths.

A.3. Computer Vision(CV) Task

We conduct experiments on the ImageNet-100, COCO2017, and ADE20K datasets using the open-

source Vanilla-VMamba-Tiny model. To ensure a fair comparison, we retrain each configuration from scratch on an 8-GPU server without employing any pre-trained weights. This avoids inconsistencies caused by potential mismatches between modified normalization layers and pre-trained parameters. Furthermore, for each dataset, we adopt the same hyperparameter settings as in the original implementation. The details are as follows:

Dataset ImageNet-100 (ima, 2019) is a curated subset of ImageNet-1k (Deng et al., 2009), comprising 100 randomly selected and semantically coherent classes. Each class contains 1,300 training images and 50 validation images, totaling 130,000 training and 5,000 validation samples. Image resolution follows that of the original ImageNet, commonly resized by cropping or scaling the short edge to 160–224 pixels. It is used to evaluate image classification performance, covering common entities such as animals, objects, and scenes, with Top-1 accuracy as the evaluation metric.

Experiment Setting The backbone consists of 14 Mamba blocks with three downsampling stages, and the layer configuration is set to [2, 2, 8, 2]. We train the model using the AdamW optimizer with a weight decay of 0.05, an initial learning rate of 5×10^{-3} , a batch size of 256, and a total of 300 epochs.

Dataset COCO 2017 (Lin et al., 2014) is one of the most widely used benchmarks for multi-task vision evaluation, featuring approximately 330,000 images, including 118,000 for training, 5,000 for validation, and 20,000 for test-dev. It includes 80 object detection categories and 91 stuff categories, with 1.5 million object instances annotated with bounding boxes. The images depict real-world scenarios with dense object layouts, occlusions, and large scale variations. It serves as a standard benchmark for object detection and instance segmentation, with mean Average Precision (mAP) used for evaluation.

Experiment Setting The backbone consists of 14 Mamba blocks with three downsampling stages, and the layer configuration is [2, 2, 9, 2]. For object detection and instance segmentation, we employ the Mask R-CNN head. The training is performed using the AdamW optimizer with a weight decay of 0.05, an initial learning rate of 1×10^{-4} , a batch size of 8, and a total of 12 epochs.

Dataset ADE-20K (Zhou et al., 2019) is a benchmark for semantic segmentation and scene parsing, aggregating over 27,000 scene images from the SUN and Places datasets. All images are annotated with pixel-level polygons, covering 150 semantic classes and over 3,000 instance-level object categories. The dataset spans a wide variety of environments, including indoor, outdoor, natural, and urban

scenes, with annotations for both visible and occluded object regions. It is the standard evaluation set for fine-grained segmentation and multi-scale understanding, using mean Intersection over Union (mIoU) as the evaluation metric.

Experiment Setting The backbone consists of 14 Mamba blocks with three downsampling stages, and the layer configuration is [2, 2, 8, 2]. For semantic segmentation, we use the UPerHead as the decoding head. Training is conducted using the AdamW optimizer with a weight decay of 0.01, an initial learning rate of 6×10^{-5} , a batch size of 32, and a total of 160,000 training iterations.

A.4. Supplementary Theoretical Backgrounds

In our section Method, we presented the definitions of the stability and optimization metrics along with the associated empirical conclusions. In this appendix, we provide supplementary theoretical discussion to further substantiate the principles underlying our analysis of training stability and optimization behavior in neural networks.

A.5. Stability Analysis

In gradient-based neural network training, instability often manifests as exploding gradients, where gradients grow excessively and lead to numerical failures. Intuitively, this phenomenon typically arises from two sources: the explosion of hidden activations (e.g., large spectral values in forward inputs or backward gradients causing NaNs in the loss), or unbounded growth in network weights due to overly large updates during backpropagation. Accordingly, our stability analysis considers both activation dynamics and weight behavior.

Due to the cumulative feature of transformations across multiple layers, the *out-put* layers in deep neural networks are particularly susceptible to numerical instabilities. Moreover, we observed a similar phenomenon in Mamba, as illustrated in Figure 11. The figure compares the activation magnitude range of the SSM outputs before and after applying normalization. It can be seen that the State Space Model (SSM) tends to amplify activation magnitudes during the forward pass, and such amplification accumulates progressively in deep networks, eventually leading to gradient explosion. Notably, in computer vision (CV) tasks, the amplification of the input \mathbf{x} by the SSM module in the absence of normalization at the `Norm2` position is significantly greater than in sequential data tasks. As a result, CV models are more susceptible to gradient explosion, as shown in Figure 12. Chiang et al. (2024) similarly observed this property of Mamba. Introducing normalization effectively mitigates this accumulation and constrains the numerical range. Therefore, we focus our analysis on the spectral properties of the forward inputs, backward gradients, and output-layer weight matrices of

out-proj layer in each Mamba block.

Previous studies have shown that spectral analysis of activations and monitoring the norms of weight matrices are effective methods for evaluating training stability in deep neural networks (DNNs) (Huang et al., 2020). Motivated by this, we adopt two representative indicators to analyze the stability of the Mamba architecture: (1) the maximum eigenvalue of the activation covariance matrix, which captures the scale and distributional dynamics of the hidden activations, and (2) the spectral norm of the output-layer weight matrices, which reflects the model’s scaling behavior and its tendency to produce unstable updates.

We denote the covariance matrix of the layer input as $\Sigma_{\mathbf{x}}^l = \mathbb{E}_{p(\mathbf{x})q(\mathbf{y}|\mathbf{x})} [\mathbf{x}^{l-1} (\mathbf{x}^{l-1})^T]$ and the covariance matrix of the layer output-gradient as $\Sigma_{\nabla \mathbf{h}}^l = \mathbb{E}_{q(\mathbf{y}|\mathbf{x})} \left[\frac{\partial \ell^T}{\partial \mathbf{h}^l} \frac{\partial \ell}{\partial \mathbf{h}^l} \right]$, where l is the l -th mamba layer.

- **Maximum eigenvalue of the input covariance matrix:**

$$\lambda_{\max}(\Sigma_{\mathbf{x}}^l) = \max \{ \lambda \in \text{Spec}(\Sigma_{\mathbf{x}}^l) \}, \quad \Sigma_{\mathbf{x}}^l = \mathbb{E}_{p(\mathbf{x})q(\mathbf{y}|\mathbf{x})} [\mathbf{x}^{l-1} (\mathbf{x}^{l-1})^T], \quad (8)$$

which measures the second-order statistics of the forward input to the layer.

- **Maximum eigenvalue of the output-gradient covariance matrix:**

$$\lambda_{\max}(\Sigma_{\nabla \mathbf{h}}^l) = \max \{ \lambda \in \text{Spec}(\Sigma_{\nabla \mathbf{h}}^l) \}, \quad \Sigma_{\nabla \mathbf{h}}^l = \mathbb{E}_{q(\mathbf{y}|\mathbf{x})} \left[\frac{\partial \ell^T}{\partial \mathbf{h}^l} \frac{\partial \ell}{\partial \mathbf{h}^l} \right], \quad (9)$$

which reflects the second-order sensitivity of the loss with respect to the layer outputs.

A.6. Optimization Analysis

In prior work, the condition number has been widely regarded as a key indicator for monitoring the optimization behavior of deep neural networks (Saarinen et al., 1993; Desjardins et al., 2015; Huang et al., 2020). It is formally defined as:

$$\kappa(\mathcal{F}_k) = \frac{\lambda_{\max}(\mathcal{F}_k)}{\lambda_{\min}(\mathcal{F}_k)}, \quad (10)$$

where \mathcal{F}_k denotes the Fisher Information Matrix (FIM) or its approximation for the k -th layer.

Intuitively, a condition number closer to 1 implies that the optimization landscape is more isotropic (i.e., closer to a spherical shape). This indicates that weight updates are more evenly distributed across the principal directions of the data, thereby facilitating more stable and efficient convergence.

FIM characterizes the curvature of the loss landscape very well (Oczkowski & Barreca, 1997; Fasina et al., 2023). One successful example is approximating the Fisher Information Matrix (FIM) of DNNs using the Kronecker-factored Approximate Curvature (K-FAC) method (Martens & Grosse, 2015). In the K-FAC approach, two assumptions are made: (1) weight gradients in different layers are assumed to be uncorrelated; (2) the input and output gradients in each layer are approximated as independent.

Under these assumptions, the full FIM can be approximated as a block diagonal matrix:

$$\mathbf{F} \approx \text{diag}(F^1, F^2, \dots, F^L),$$

where F^l is the sub-FIM corresponding to the parameters in the l -th layer, computed as:

$$F^l = \mathbb{E}_{p(\mathbf{x}), q(\mathbf{y}|\mathbf{x})} \left(\left(\mathbf{x}^{l-1} (\mathbf{x}^{l-1})^T \right) \otimes \left(\frac{\partial \ell^T}{\partial \mathbf{h}^l} \frac{\partial \ell}{\partial \mathbf{h}^l} \right) \right) \\ \approx \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} \left[\mathbf{x}^{l-1} (\mathbf{x}^{l-1})^T \right] \otimes \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim p(\mathbf{x})q(\mathbf{y}|\mathbf{x})} \left[\frac{\partial \ell^T}{\partial \mathbf{h}^l} \frac{\partial \ell}{\partial \mathbf{h}^l} \right]$$

Properties of the Kronecker-structured FIM. Given the Kronecker-factored approximation of the Fisher Information Matrix (FIM) for the l -th layer as

$$F^l = A \otimes B,$$

where A and B are symmetric positive semi-definite matrices representing the input and output-gradient statistics respectively, the spectral properties of F^l satisfy the following:

- **Eigenvalues.** The eigenvalues of F^l are the pairwise products of the eigenvalues of A and B :

$$\lambda(F^l) = \{\lambda_i(A) \cdot \lambda_j(B) \mid \lambda_i \in \lambda(A), \lambda_j \in \lambda(B)\}.$$

In particular, the maximum eigenvalue satisfies:

$$\lambda_{\max}(F^l) = \lambda_{\max}(A) \cdot \lambda_{\max}(B).$$

- **Condition Number.** The condition number of F^l equals the product of the condition numbers of A and B :

$$\kappa(F^l) = \kappa(A \otimes B) = \kappa(A) \cdot \kappa(B).$$

B. Empirical Evidence for the Role of Norm2 in Enhancing Training Stability

Previous studies have shown that Layer Normalization (LN) effectively stabilizes model convergence (Ba, 2016; Xu et al., 2019; Zhang & Sennrich, 2019). In particular, research within Transformer architectures has demonstrated

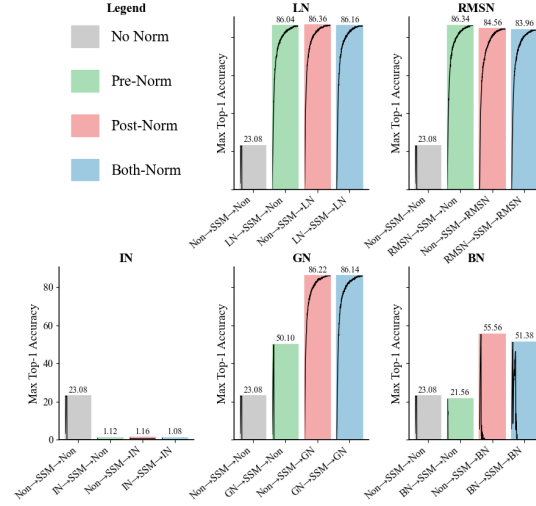


Figure 10. Stability comparison across different normalization methods applied to the SSM block.

that applying LN before the attention module can yield superior performance (Kim et al., 2025; Xiong et al., 2020; Shleifer et al., 2022). This raises the question: does a similar principle hold for the Mamba architecture? Specifically, is LN more effective when placed before the SSM block (Norm1), after the SSM block (Norm2), or on both sides?

To investigate this, we conducted experiments on the ImageNet-100 dataset with three configurations: applying LN before the SSM module (Norm1), after the SSM module (Norm2), and on both sides (Norm1 and Norm2). All experiments were conducted using the same model architecture and training settings to ensure fairness. The results are presented in Figure 13.

Key observations from Figure 13:

1. **Normalization can stabilize model training:** The baseline configuration (Non-SSM-Non) fails to converge, achieving only 23.08% Top-1 accuracy. This indicates that without normalization, the amplification of activations and gradient explosion in the SSM severely disrupt training. In contrast, all configurations with normalization successfully converge and achieve Top-1 accuracy above 86%.
2. **Post-SSM LN (Norm2) yields the best performance:** When LN is applied only after the SSM (Non-SSM-LN), the model achieves the highest Top-1 accuracy of **86.36%** and exhibits the steepest convergence curve. This suggests that Norm2 directly suppresses the activation norm explosion from the SSM, significantly improving training stability and conver-

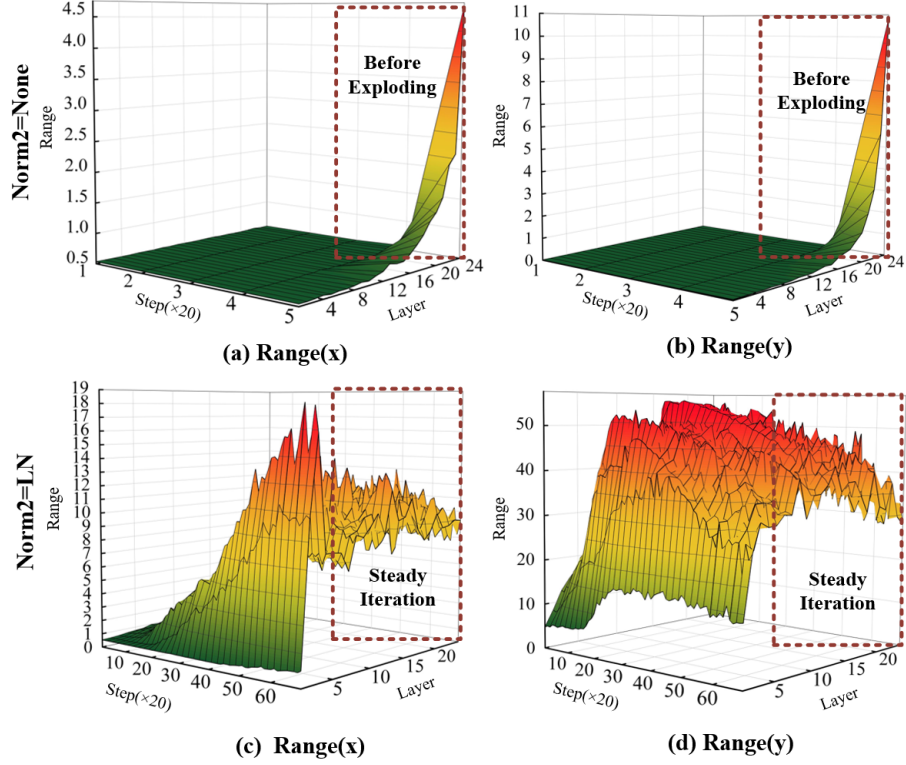


Figure 11. Comparison of the input (x) and output (y) range of the SSM module on the WIKITEXT-103 dataset. The x -axis denotes training steps (logged every 20 steps), and the y -axis indicates the layer index of the Mamba block. Subfigures (a) and (b) show the range evolution when no normalization is applied after the SSM layer (Norm2=None), while (c) and (d) present the results with LayerNorm applied at Norm2. Applying LN leads to steady iteration dynamics across layers, whereas the absence of normalization results in unstable growth before divergence.

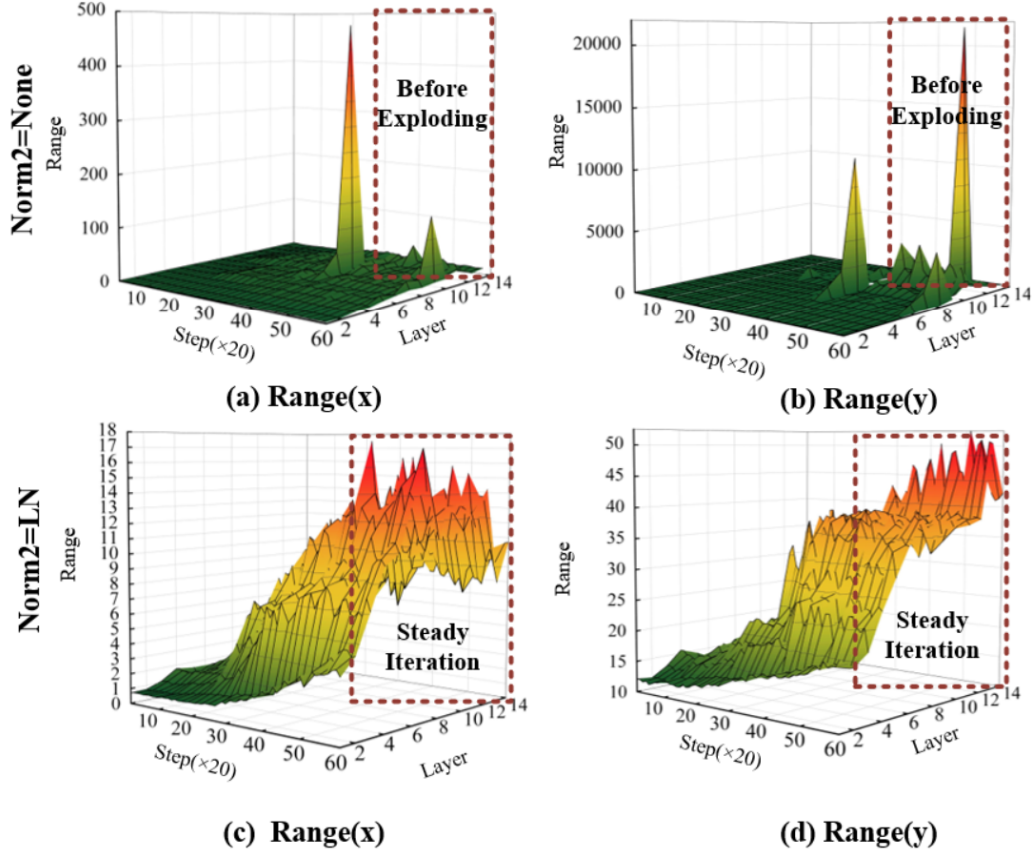


Figure 12. Comparison of the input (x) and output (y) range of the SSM module on the IMAGENET dataset. The x -axis denotes training steps (logged every 20 steps), and the y -axis indicates the layer index of the Mamba block. Subfigures (a) and (b) show the range evolution when no normalization is applied after the SSM layer (Norm2=None), while (c) and (d) present the results with LayerNorm applied at Norm2. Applying LN leads to steady iteration dynamics across layers, whereas the absence of normalization results in unstable growth before divergence.

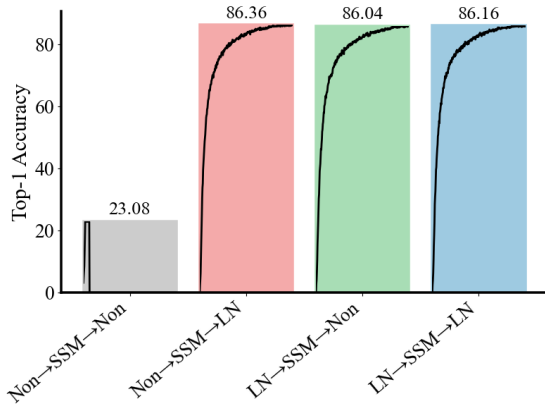


Figure 13. Comparison of training stability with LN applied at different positions in the SSM block. Each bar represents the final Top-1 accuracy, and the inner curve shows the accuracy evolution during training.

gence speed.

3. **Pre-SSM LN (Norm1) is slightly less effective:** Applying LN only before the SSM (LN→SSM→Non) results in a slightly lower peak accuracy of **86.04%** and slower convergence. This implies that Norm1 helps improve the numerical condition of the input features but is less effective than post-normalization in stabilizing the network.
4. **Dual LN provides stable convergence but no additional benefit:** Using LN on both sides of the SSM (LN→SSM→LN) yields a final accuracy of **86.16%**, which is between the performances of Norm1 and Norm2 alone. The convergence speed is also intermediate. This suggests that dual normalization does not provide additive benefits and may slightly impair feature expressiveness due to over-normalization.

In summary, LN is effective in stabilizing training in Mamba-based models. However, unlike Transformer models where pre-normalization is often optimal, **post-normalization (Norm2) after the SSM block achieves better stability and performance in Mamba**. While pre-normalization (Norm1) accelerates convergence, it is slightly less effective, and dual normalization brings no additional gains.

B.1. Comparison with Other Normalization Methods

To further validate the generality of post-normalization effectiveness beyond LN, we also evaluated other commonly used normalization strategies, including RMSNorm and GroupNorm. The experimental results are shown in Figure 10.

Key observations from Figure 10:

- **RMSNorm and LayerNorm:** Both normalization methods are capable of stabilizing training regardless of placement (pre or post-SSM). For GroupNorm, however, only post-normalization configurations (Non→SSM→GN and GN→SSM→GN) result in stable convergence with Top-1 accuracy exceeding **86.0%**. The pre-normalization setting (GN→SSM→Non) diverges after several training steps, reaching only **50.1%** accuracy.
- **InstanceNorm and BatchNorm:** Neither method achieves convergence under any configuration. Post-normalization variants (Non→SSM→IN and Non→SSM→BN) initially improve performance but subsequently diverge, with Top-1 accuracy peaking at only **1.12%** and **21.56%**, respectively, before collapsing.

These results reinforce the finding that, unlike in Transformer architectures, **post-SSM normalization (Norm2) is particularly effective for ensuring training stability in Mamba models**, especially when using LayerNorm. Accordingly, we adopt $\text{Norm2} = \text{LN}$ as the default normalization configuration in our main experiments.

C. Empirical Evidence for the Role of Norm1 in Improving Optimization Condition

C.1. Effect of Batch Normalization on Accelerating Convergence and Comparison with Other Normalization Techniques

Having established a stable training setup with post-SSM Layer Normalization (Norm2 = LN), we further investigate the effect of Batch Normalization (BN) on accelerating convergence. Specifically, we use the Non→SSM→LN configuration as our baseline, and introduce BN before the SSM block (BN→SSM→LN, i.e., Norm1 = BN). The experimental results are illustrated in Figure 14.

Observations from Figure 14:

With Norm2 fixed as LN to ensure stable training, introducing BN at Norm1 (red line) significantly accelerates convergence compared to the baseline Non+LN (blue line):

- **Faster early-stage convergence:**
 - At epoch 10, BN+LN achieves a Top-1 accuracy of approximately 35.6%, compared to 34.8% for Non+LN.
 - At epoch 20, BN+LN reaches about 62.3%, approximately 4 percentage points higher than the baseline.

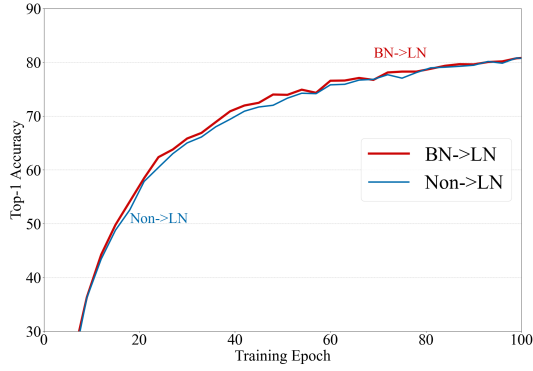


Figure 14. Convergence speed comparison between BN+LN and Non+LN configurations.

- **Mid-stage acceleration:**

- From epoch 20 to 50, BN+LN maintains a steeper ascent, surpassing 70% accuracy earlier than the baseline.
- By epoch 60, the accuracy curve of BN+LN stabilizes around 75%, with reduced oscillation.

These findings suggest that BN→SSM→LN enables larger gradient steps during the early training phase, allowing the model to enter the high-accuracy regime faster and resulting in significantly improved convergence speed over Non→SSM→LN.

C.1.1.1. COMPARISON WITH OTHER NORM1 CONFIGURATIONS

To further assess the acceleration effect across different normalization strategies, we replaced Norm1 with other commonly used normalization methods, including Group Normalization (GN), Root Mean Square Normalization (RMSN), Instance Normalization (IN), and Layer Normalization (LN). We compared their Top-1 accuracy and the steepness of the convergence curves. The results are presented in Figure 15.

Key findings from Figure 15:

- **BN+LN achieves the best performance:**

- Final Top-1 accuracy reaches **86.72%**, surpassing the baseline Non+LN (86.36%).
- The learning curve exhibits the steepest slope in the 0–80% range, indicating the most significant early acceleration.

- **Other combinations:**

- GN→SSM→LN and LN→SSM→LN also provide faster convergence than the baseline, but with gentler slopes during the first 20–30 epochs.

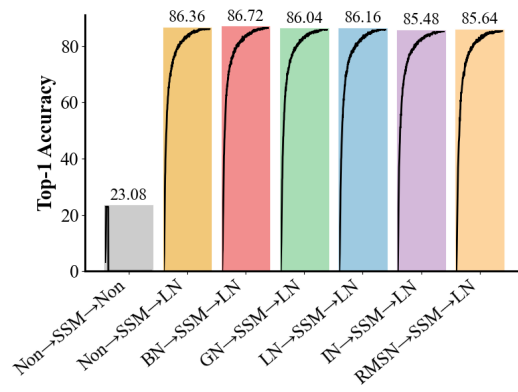


Figure 15. Convergence speed comparison with various Norm1 methods combined with Norm2 = LN.

- IN→SSM→LN and RMSN→SSM→LN exhibit the flattest early-stage slopes and the slowest overall convergence, with the lowest final accuracy.

In summary, among all evaluated Norm1 strategies, the **BatchNorm** → **SSM** → **LayerNorm** configuration not only maintains high final accuracy but also maximizes convergence speed, making it the most effective combination for efficient training.