

AN EMPIRICAL STUDY ON NORMALIZATION IN MAMBA

Anonymous authors

Paper under double-blind review

ABSTRACT

Normalization layers are crucial for enhancing both the training efficiency and stability of deep neural network architectures. The recently proposed Mamba network has shown considerable promise in achieving competitive results alongside Transformers. However, ensuring training stability in Mamba, as in many deep architectures, remains a significant challenge, where normalization techniques are crucial for solving this issue. In this paper, we conduct a systematic investigation into the effects of various normalization techniques and their combinations within the Mamba architecture. Our analysis encompasses both long sequence modeling and image classification tasks. For long sequence modeling, we perform extensive experiments to assess the impact of applying normalization layers both before and after S6 Modules. The results show that normalization across layers leads to enhanced training stability and improved model performance. Furthermore, we validate these findings on large-scale V-Mamba models and offer practical suggestions for selecting appropriate normalization techniques. We hope that our insights will contribute to mitigating training instabilities in deep learning and fostering the development of more robust architectures. All code and models used in this study will be open-sourced on GitHub.

1 INTRODUCTION

The Mamba architecture (Gu & Dao, 2023) is a state-of-the-art model designed to efficiently handle long sequences in sequence modeling (Gu & Dao, 2023), leveraging structured state-space models (SSM) (Gu et al., 2021). However, training Mamba presents significant challenges due to the inherent unpredictability of state space model (SSM). Normalization plays a crucial role in the effective training of deep neural networks (Huang et al., 2023), particularly for architectures like Mamba. There are various usage of normalization techniques in the variants of Mamba architecture (Liu et al., 2024; Liang et al., 2024), but without providing the sufficient evidences for supporting why one should use like that. Therefore, it is essential to conduct a comprehensive investigation into various normalization strategies to optimize Mamba’s performance. In this study, we examine the impact of several widely used normalization techniques, including Batch Normalization (BN) (Ioffe, 2015), Layer Normalization (LN) (Ba, 2016), Group Normalization (GN) (Wu & He, 2018), Instance Normalization (IN) (Huang & Belongie, 2017), and Root Mean Square Normalization (RMSN) (Zhang & Sennrich, 2019). We systematically evaluate the effect of normalization type, position, and combinations on Mamba’s performance and stability in sequence modeling and image classification tasks.

For sequence modeling tasks, we found that placing normalization before the S6 layer led to performance degradation while placing it after the S6 Module led to performance improvement. Certain combinations of normalization techniques applied before and after the S6 Module enhance the model’s performance, while some can be detrimental. Our analysis, rooted in the scale-invariant property of deep learning (Papyan, 2018), revealed that the weight matrices (Huang et al., 2020) in deeper Mamba blocks exhibited significantly larger L2 norms (Luo et al., 2016) than earlier layers. Introducing normalization before the S6 Module exacerbated this issue, whereas placing normalization after the S6 Module helped the model enter a state of *Weight Domination*, which is beneficial to keeping scale-invariant property. However, *Weight Domination* also limited the extent to which weights could be updated. By applying normalization before the S6 Module, we found that this *Weight Domination* phenomenon could be alleviated, leading to more effective parame-

054 ter updates under stable training conditions. Finally, we propose an intution for selecting optimal
 055 normalization combinations to improve both the efficiency and stability of Mamba.
 056

057 For image classification tasks, we validated the conclusions derived from sequence modeling. We
 058 explored the impact of position on performance by examining BN, IN, LN, RMSN, GN, only RMSN
 059 underperformed. However, by positioning the normalization layers before and after the Mamba
 060 blocks, we aimed to determine whether the placement of these techniques affected network per-
 061 formance. For IN, GN, and LN, placing normalization before or after the Mamba blocks consistently
 062 reduced validation loss and improved accuracy. However, when BN or RMSN were placed before
 063 the Mamba blocks, validation loss increased and accuracy declined, while placing these layers after
 064 the Mamba blocks resulted in unstable performance. These results suggest that LN normalization
 065 techniques are robust to positional changes, whereas BN and RMSN are sensitive to position.
 066

067 Finally, motivated by those findings in sequence modeling, we investigated combinations of nor-
 068 malization techniques before and after the Mamba blocks to maximize network performance. We
 069 explored extensive combinations, such as BN+RMSN, LN+BN, LN+RMSN, etc. Surprisingly, even
 070 normalizations that previously detrimented performance, such as RMSN and BN led to combined
 071 improvements. Notably, the RMSN+BN combination yielded faster and more significant improve-
 072 ments than the other combinations like LN+LN. Compared to the self-combination of RMSN, IN,
 073 GN, and LN, cross-combined with BN resulted in faster and greater reductions in validation loss and
 074 higher accuracy. However, LN+RMSN underperformed relative to LN+LN. These results demon-
 075 strate: (a) the order of normalization placement affects the Mamba network’s performance; (b)
 076 placing LN before the Mamba block and another normalization technique after helps mitigate the
 077 negative effects of certain normalizations; (c) BN is particularly effective when placed after the
 078 Mamba block and is an essential component for enhancing the Mamba network’s performance when
 079 combined with other normalization techniques.
 080

Contributions In this work, we systematically investigate widely used normalization techniques
 081 and aim to address the following questions:
 082

(1) Normalization Position: Where should normalization be placed in the model to ensure high
 083 performance and stable training? We found that applying normalization layer after the S6 Module
 084 SSM yields better results(if only one can be choosen). However, combining certain normalization
 085 types can achieve more outstanding performance.
 086

(2) Normalization Combanition: What are the effects of different normalization combinations? And
 087 what is the intuition behind these effects? We found that certain combinations of different normal-
 088 ization techniques produce excellent results.
 089

(3) Combination intution: How can we choose the proper combination?We propose an intution
 090 for harmonizing normalization strategies in deep architectures, providing practical guidelines for
 091 selecting normalization methods.
 092

2 PRELIMINARY

2.1 MAMBA FOR SEQUENCE MODELING

096 Mamba (Gu & Dao, 2023) is a state-of-the-art architecture based on the State Space Model
 097 (SSM) (Hamilton, 1994). For one-dimensional sequence modeling, we define the sequence as
 098 $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_L\} \in \mathbb{R}^{L \times D}$, where L represents the sequence length and D denotes the number
 099 of feature channels. Each element $\mathbf{x}_l = \{x_l^1, x_l^2, \dots, x_l^D\} \in \mathbb{R}^D$ contains D feature values at time
 100 step l . For a mini-batch of data with m samples, the input tensor is typically shaped as (m, D, L) .
 101

Not only does Mamba excel in handling one-dimensional sequences (Gu & Dao, 2023), but it also
 102 extends effectively to two-dimensional image data (Liu et al., 2024).
 103

Because image data usually undergoes a flatten operation to transform the $h * w$ dimension into
 104 sequence length L before entering the model. In particular, to allow each pixel to better gather
 105 information from all directions and establish a global receptive field in the 2D space, a 2D-Selective-
 106 Scan (SS2D) method is proposed in VMamba. In this method, images are divided into patches,
 107 and each patch is processed as a sequence along four different traversal paths (Cross-Scan) using
 108 parallel S6 blocks (Mamba). The processed sequences are then reshaped and merged to form the
 109

108 output map (Cross-Merge), thus enabling Mamba to efficiently process both 1D and 2D data. In
 109 our experiments, to better study the Mamba block’s nature, we removed the FFN block, which is an
 110 MLP branch following the VSS Block in VMamba. For VMamba details, please refer to (Liu et al.,
 111 2024).

112

113 2.2 NORMALIZATION

114

115 Normalization is a crucial technique in deep learning (Hinton & Salakhutdinov, 2006; Ioffe, 2015;
 116 He et al., 2015; Huang et al., 2023) that aims to stabilize and accelerate the training process of
 117 neural networks (Desjardins et al., 2015), as defined in Equation 1. It helps adjusting input data
 118 distribution or intermediate activations into normal distribution (Huang et al., 2023) that is beneficial
 119 to the convergence of model training. Furthermore, normalization methods, such as BN, LN and GN
 120 address issues like internal covariate shift (Wang et al., 2022), which occurs when the distribution
 121 of inputs to each layer changes during training (Wang et al., 2022). Models can maintain consistent
 122 representation learning and effectively manage varying input scales by incorporating normalization
 123 layers. Thus, selecting the appropriate normalization method according to specific requirements is
 124 essential for training with complex data or large parameters.

125

$$\text{Norm}(x) = \frac{x - E[x]}{\sqrt{Var[x] + \epsilon}} \cdot \gamma + \beta \quad (1)$$

126

127 Various normalization variants have been proposed, including Iterative Normalization (ItN) for
 128 addressing normalization computational inefficiencies (Huang et al., 2019), Scale Normalization
 129 (Scale) for ensuring consistent value scaling of features (Papyan, 2018). However, the five most
 130 widely adopted methods by researchers are still BN, LN, GN, IN, and RMSN, as they have been
 131 well-integrated into pytorch libraries and cater for most models. Accordingly, our experiments fo-
 132 cus on these commonly used normalization methods.

133

134 Huang et al. (2023) proposed a unified taxonomy for understanding the similarities and differences
 135 between these approaches, specifically, Normalization Representation Area Partitioning (NAP),
 136 Normalization Operation (NOP), and Normalization Representation Recovery (NRR). The opera-
 137 tional details of these mainstream techniques are clearly presented in Table 1.

138

139 Table 1: Details of Mainstream Normalization techniques: We take sequence $\mathbf{X} \in \mathbb{R}^{m \times D \times L}$ as an
 140 example. The NAP operation determines how \mathbf{X} is reshaped into $\mathbf{X} \in \mathbb{R}^{S_1 \times S_2}$, where S_2 indexes
 141 the set of samples used to compute the statistics. For example, $\mathbf{X} = \Pi_{BN}(\mathbf{X}) \in \mathbb{R}^{D \times (mL)}$, where
 142 mL indicates that the statistics (mean and variance) are calculated along the batch, sequence length
 143 (time steps)

144

Method	NAP	NOP	NRR
Batch Normalization (BN)	$\Pi(\mathbf{X}) \in \mathbb{R}^{D \times mL}$	Standardizing	Learnable $\gamma, \beta \in \mathbb{R}^D$
Instance Normalization (IN)	$\Pi(\mathbf{X}) \in \mathbb{R}^{mD \times L}$	Standardizing	Learnable $\gamma, \beta \in \mathbb{R}^D$
Group Normalization (GN)	$\Pi(\mathbf{X}) \in \mathbb{R}^{mg_D \times s_DL}$	Standardizing	Learnable $\gamma, \beta \in \mathbb{R}^D$
Layer Normalization (LN)	$\Pi(\mathbf{X}) \in \mathbb{R}^{m \times DL}$	Standardizing	Learnable $\gamma, \beta \in \mathbb{R}^D$
Root Mean Square Norm (RMSN)	$\Pi(\mathbf{X}) \in \mathbb{R}^{m \times DL}$	Scaling	Learnable $\gamma \in \mathbb{R}^D$

153

154

155

156 3 THE MAMBA LANDSCAPE

157

158 Due to the inherent unpredictability of state space transitions, especially when introducing the multi-
 159 head mechanism, Mamba is challenging to train stably during the training process. To mitigate this
 160 issue, it becomes essential to implement effective strategies for stabilizing the training dynamics.
 161 One such key solution involves incorporating an extra normalization layers in the architecture. In
 Section 3.1, we describe the process of sequence modeling using Mamba architecture and specify

the positions of the normalization layers we investigated (In Figure 1, the orange layers: Normalization1, before inprojection layer which norms the input activations and Normalization2 ,before outprojection layer which norms the SSM output.). Followed by a detailed discussion of the normalization techniques explored in Section 3.2.

3.1 MAMBA LANDSCAPE DETAILS

As shown in Figure 1, the input data is first processed in the data factory, the sequence data will be extracted with spatio-temporal features, and the image data will undergo scanning and flatten operations to become sequence data. After being embedded, the data will be projected into high-dimensional feature space through the inprojection layer. A convolution operation is then applied to fuse information from adjacent time steps, followed by processing through the SSM layer and variants including S4 (Gu et al., 2021), S6 (Gu & Dao, 2023), SSD (Dao & Gu, 2024). The output is subsequently combined with the data from the feedforward MLP channel, undergoes an extra normalization layer, and is finally projected to the desired dimensions through the outprojection layer.

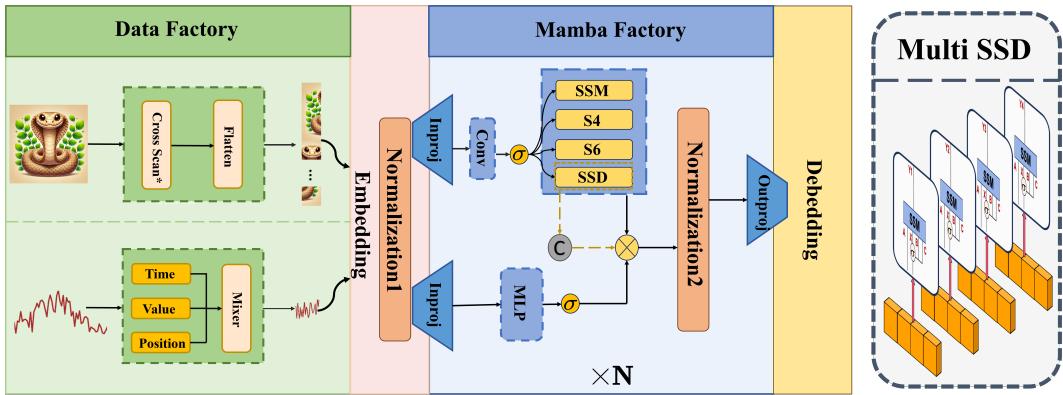


Figure 1: Illustration Data Flow in Mamba Architecture and Our Research Focus:

In previous Mamba variants, many of them has introduced the normalization operation right after the S6 Module. But there is no unified study on where and what normalization to introduce in the architecture. We have extensively explored the the type, position, combination of different normalization techniques.

3.2 EXTRA NORMALIZATION

The normalization in Mamba can be broadly categorized into two types:

1. Normalization applied to the input activations before entering the S6 Module.
2. Normalization applied to the output activations after the S6 Module.

We first define the notation of normalization to describe the position of normalization layers within a block. Let Norm1→Norm2 denote the normalization configuration within a block, where Norm1 refers to the normalization applied before inprojection layer, and Norm2 refers to the normalization applied before the output projection layer. For instance, the notation BN→IN represents the application of Batch Normalization before the inproj and Instance Normalization before the outproj, as depicted in Figure 2.

3.3 DATASETS

For long sequence modeling, we use the LRA ListOps task (Tay et al., 2021), a challenging benchmark designed to test a model’s capacity to handle long-range dependencies. The task requires the model to classify sequences of nested operations applied to integers, making it a strong indicator

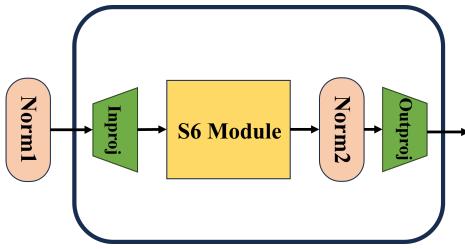


Figure 2: Norm Location: Norm1 is responsible for normalizing the input activations, while Norm2 normalizes the output of the S6 Module (Mamba Block).

of a model’s ability to process long sequences. Furthermore, we extend the difficulty by utilizing a longer version of ListOps with an input sequence length of 2000 tokens, making the task more demanding for the model to capture distant dependencies effectively.

For vision tasks, we use the ImageNet-100 (Krizhevsky et al., 2012) classification dataset. ImageNet-100 is a randomly selected subset from the ImageNet-1k dataset of the 2012 Large Scale Visual Recognition Challenge. It contains 100 categories, spanning various objects and scenes, ensuring diversity in vision tasks. The training set has 1,300 images per category, and the validation set contains 50 images per category, totaling 135k images.

3.4 EXPLORING THE MAMBA NORMALIZATION LANDSCAPE

We conduct a comprehensive investigation of commonly used normalization techniques, analyzing their type, position and combinations in both long sequence modeling and image classification tasks. In Section 4.1, we fix everything but combination evaluate the effects of applying normalization before or after the S6 Module. Our results indicate that applying normalization after the S6 Module usually improves model , while appending it before the S6 Module leads to degraded performance. In Section 4.2, we examine the impact of combining different normalization techniques, with certain combinations resulting in significant performance improvements. A detailed analysis of these findings will follow. Finally, in Section 4.3, we validate these results using the Mamba architecture for image classification, further confirming the observed trends in sequence modeling.

4 EXPERIMENTAL RESULTS AND DISCUSSION

4.1 NORMALIZATION POSITION

Table 2: Accuracy Comparison: The left column shows normalization applied before the S6 Module, the middle column shows it applied after, and the right column shows normalization applied both before and after the S6 Module.Underlined None → None is the baseline model.The first arrow(\uparrow) compares with the first column, and the second arrow($\uparrow\downarrow$) compares with the second column.

Normalization	Accuracy	Normalization	Accuracy	Normalization	Accuracy
<u>None</u> → None	0.390	<u>None</u> → None	0.390	<u>None</u> → None	0.390
BN → None	0.389	None → BN	0.426(\uparrow)	BN → BN	0.439($\uparrow\downarrow$)
IN → None	0.357	None → IN	0.380(\uparrow)	IN → IN	0.406($\uparrow\downarrow$)
LN → None	0.374	None → LN	0.394(\uparrow)	LN → LN	0.381 ($\uparrow\downarrow$)
RMSN → None	0.379	None → RMSN	0.394(\uparrow)	RMSN → RMSN	0.401($\uparrow\downarrow$)
GN → None	0.390	None → GN	0.394(\uparrow)	GN → GN	0.414($\uparrow\downarrow$)

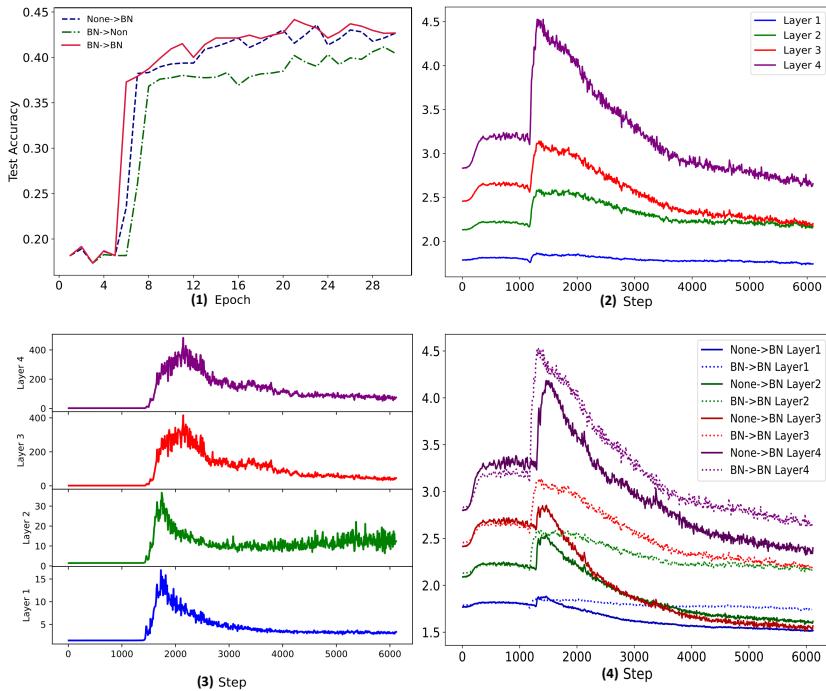
The results of experiments are summarized in Table 2. We can draw three conclusions:

1. We observed that adding normalization before the S6 Module weakens the model’s predictive performance, with Instance Normalization (IN) showing the most significant decline.

270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
2. However, when normalization is applied after the S6 Module, the model's predictive accuracy improves compared to normalization applied before the S6 Module, and all results outperform those without normalization.

3. When the same normalization is applied both before and after the S6 Module, there is a further improvement in predictive accuracy. Except for Layer Normalization (LN), all other normalization techniques show additional improvements when applied both before and after the S6 Module.

We plotted the L2 norms of the weights in each Mamba block and conducted an in-depth analysis of the effects of different normalization techniques in a 4-layer model in the following Figure 3. We



303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
Figure 3: The Performance and L2 norm about BN→BN compared with None→BN and BN→None

present the test accuracy for three configurations: None→BN, BN→None, and BN→BN. Subfigure in subfigure (1), show the L2 norms of the weight parameters for each block in the None→BN configuration in subfigure (2) , depict the L2 norms of the weight parameters for each block in the BN→None configuration in subfigure (3) and illustrate the L2 norms of the weight parameters for each block in the BN→BN configuration in subfigure (4). The solid line represents the results of the None→BN configuration, while the dashed line corresponds to the results of the BN→BN configuration.

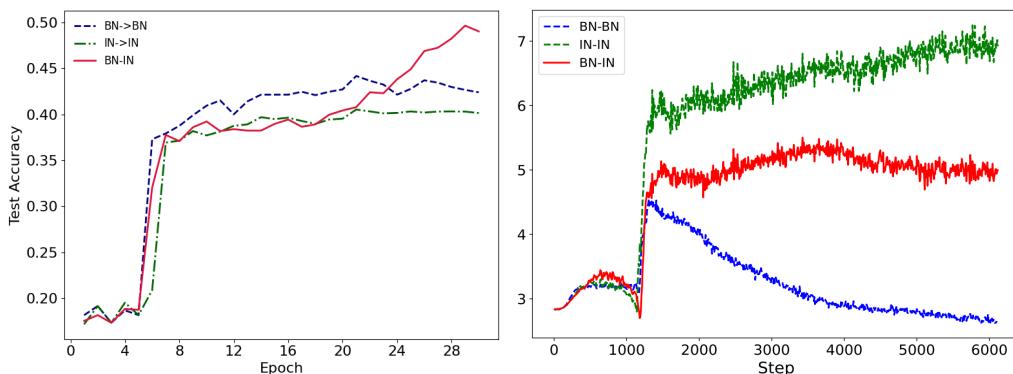
We observe that in the BN→None configuration, as the network layer becomes deeper, there is a significant disparity in the L2 norms of the weight parameters across layers, indicating a large-scale variation between layers. This suggests that the later layers of Block dominate the weighting, which is not conducive to network stability. In contrast, in the None→BN configuration, the magnitude of the L2 norms of the weight parameters remain within the same across layers. This indicates smaller scale variations between layers, leading to more stable training. Last but not least, we also observe that the updates in the weight norms are relatively limited across all layers, suggesting that the model enters a state of *Weight Domination*, making it harder to update the parameters effectively. In the BN→BN configuration, the weight norms across layers remain within the same magnitude but with larger magnitude of the update in the L2 norms. This implies that applying normalization before the S6 Module mitigates the issue of *Weight Domination*, facilitating more effective parameter updates. For more cases, please refer to Appendix C

324 4.2 NORMALIZATION COMBINATION
325

326 Additionally, we tried different normalization combinations and listed the excelled ones in the ta-
 327 ble3.Then,we provide valuable guidance for selecting optimal normalization combinations. For
 328 instance, here is the behavior of BN and IN combination. Our analysis reveals that these two nor-
 329 malization techniques exhibit complementary effects on model weights in the final block,which we
 330 refer to as "harmonic structure".We observed that the weight matrix Norm updates in different di-
 331 rections and has a large spacing when the two normalizations act alone. The Norm of BN→IN, on
 332 the other hand, is exactly the balance of the two. This leads to a 10% improvement in performance
 333 compared to using either normalization individually (see Left Figure 4).

334 Table 3: Combination Comparison: Here are the top 5 combinations and the full results can be found
 335 in the table7.
 336

337	Normalization	Accuracy
339	BN → IN	0.469
340	IN → GN	0.468
341	GN → IN	0.448
342	RMSN → BN	0.448
343	BN → BN	0.439
344	<u>None</u> → None	0.390



359 Figure 4: The Performance and L2 in layer 4 of BN→IN
 360
 361

362 4.3 VISION MAMBA
363

364 For vision tasks, we use a tiny version of the VMamba model for image classification . The model
 365 we use has a depth of 15 Mamba blocks and contains three downsampling layers. The original
 366 Mamba blocks consist of normalization layers, inprojection layer, convolutional layers, activation
 367 layers, the S6 module, and feedforward MLP layers, with the same normalization layer applied both
 368 before and after the S6 module. Due to the multiple components involved, different combinations
 369 can form various styles of Mamba blocks.

370 Given the important role of normalization layers in network design and training, we first investigate
 371 the impact of the normalization layer's position within the Mamba block on model performance.
 372 Since the core of the Mamba block is the S6 module, we consider the effect of the normalization
 373 layer's position relative to the S6 module, i.e., whether it is placed before or after the S6 module.
 374 Currently, the same normalization layer is used both before and after the S6 module, leading to a
 375 rather uniform approach. This prompted us to explore whether other diverse combinations of nor-
 376 malization layers could improve model performance, and we further experimented with the effects
 377 of combining different normalization layers. The experimental results are shown in Table 3 and
 Table 4, respectively.

378
 379 Table 4: The impact of normalization layer positioning relative to S6 module, evaluated with five
 380 representative normalization layers. LN → LN is the baseline model adopted in original Vmamba.
 381 Accuracy1 represents top 1 accuracy. The first arrow (\uparrow) compares with the first column, and the
 382 second arrow ($\uparrow\downarrow$) compares with the second column.

Normalization	Accuracy	Normalization	Accuracy	Normalization	Accuracy
None → None	0.196	None → None	0.196	None → None	0.196
BN → None	0.167	None → BN	0.678(\uparrow)	BN → BN	0.746($\uparrow\downarrow$)
IN → None	0.702	None → IN	0.838(\uparrow)	IN → IN	0.863($\uparrow\downarrow$)
LN → None	0.864	None → LN	0.867(\uparrow)	LN → LN	0.866
RMSN → None	0.294	None → RMSN	0.176(\downarrow)	RMSN → RMSN	0.072($\downarrow\downarrow$)
GN → None	0.660	None → GN	0.868(\uparrow)	GN → GN	0.863($\uparrow\downarrow$)

391 From Table 4, we found that placing the normalization layer after the S6 module significantly im-
 392 proves the classification performance of the VMamba model. Among them, GN and LN are superior
 393 choices. When they are placed after the S6 module, the model achieved accuracies of 0.868 and
 394 0.867, respectively, while IN also performed well with 0.838. Conversely, placing the normaliza-
 395 tion layer before the S6 module led to a significant performance drop. For example, BN → None
 396 only achieved 0.167, whereas None → BN reached 0.678. RMSN showed the worst performance,
 397 regardless of whether it was placed before or after SS, with lower accuracy compared to other meth-
 398 ods. These two conclusions corroborate what we drawn in the sequence experiments.

399 Table 5: The impact of different normalization layer combinations on the VMamba model, sorted
 400 by accuracy.

Normalization	Accuracy1	Normalization	Accuracy1	Normalization	Accuracy1
RMSN → BN	0.876	BN → RMSN	0.845	None → BN	0.678
GN → BN	0.872	IN → IN	0.837	RMSN → None	0.294
LN → BN	0.871	GN → RMSN	0.836	None → None	0.196
IN → BN	0.867	LN → RMSN	0.836	None → RMSN	0.176
<u>LN → LN</u>	<u>0.866</u>	IN → RMSN	0.821	BN → None	0.167
GN → GN	0.863	BN → BN	0.746	RMSN → RMSN	0.072

410 Table 5 results show that combining different normalization layers has a significant impact on the
 411 performance of the S6 module in the VMamba model for image classification tasks. Notably, the
 412 combination of RMSN → BN achieved the highest accuracy of 0.876, while standalone RMSN →
 413 RMSN performed the worst, only 0.072. This indicates that a proper combination of normalization
 414 layers can greatly enhance model performance.

415 We believe that the BN normalization layer is better suited for handling the features processed by
 416 the S6 module, thus improving the overall model performance. In contrast, using RMSN alone may
 417 disrupt the input feature distribution, leading to a significant performance drop.

418 We also explored the impact of normalization at different positions on training the VMamba model.
 419 The training loss and validation accuracy are shown in Figure 6. In Figure 6, the same color rep-
 420 presents the same normalization layer, with dashed and solid lines indicating normalization applied
 421 before and after the S6 module, respectively. From 6, we found that shifting the position of GN
 422 and IN normalization from before to after the S6 module transforms their negative impact on train-
 423 ing acceleration into a positive one. Additionally, LN normalization shows robustness to positional
 424 changes, as it consistently accelerates model convergence whether applied before or after the S6
 425 module.

426 To more intuitively demonstrate the impact of different normalization layers on training VMamba,
 427 we also plotted the curves of training loss and validation accuracy during the training of VMamba
 428 models with different normalization layers and positions, as shown in Figure 5. BN and RMS normalization
 429 layers degraded the performance of the Mamba network architecture, while LN, GN, and IN improved its performance. Additionally, we plotted the training loss and validation accu-
 430 racy when combining the performance-degrading BN and RMSN layers with other normalization
 431 layers, as shown in Figure 7. We observed that their negative effects were eliminated when com-

bined with other normalization layers. Furthermore, the normalization layers combined with BN all outperformed their individual performance.

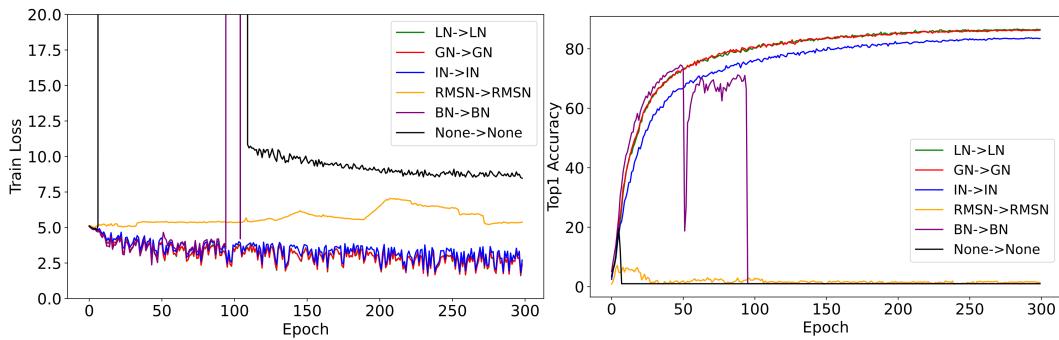


Figure 5: Impact of different Normalization on Vmamba performance. BN and RMS normalization layers degraded the performance, while LN, GN, and IN improved VMamba performance.

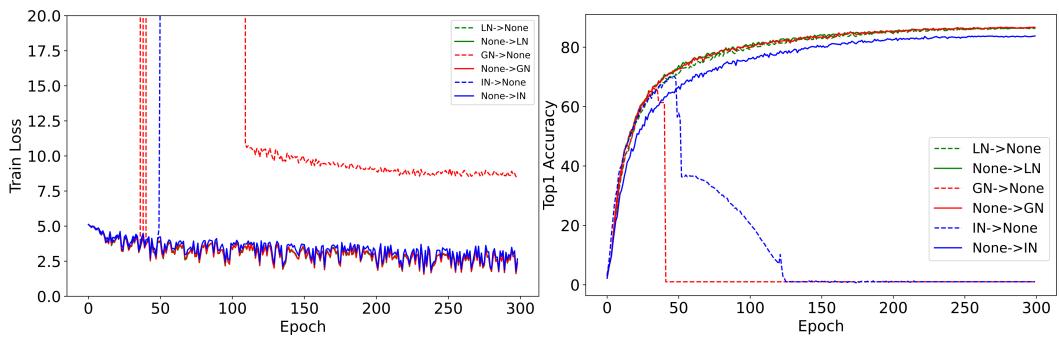


Figure 6: Influence of Location of Normalization: The solid line represents the results where normalization is applied after the S6 Module, while the dashed line represents the results where normalization is applied before the S6 Module

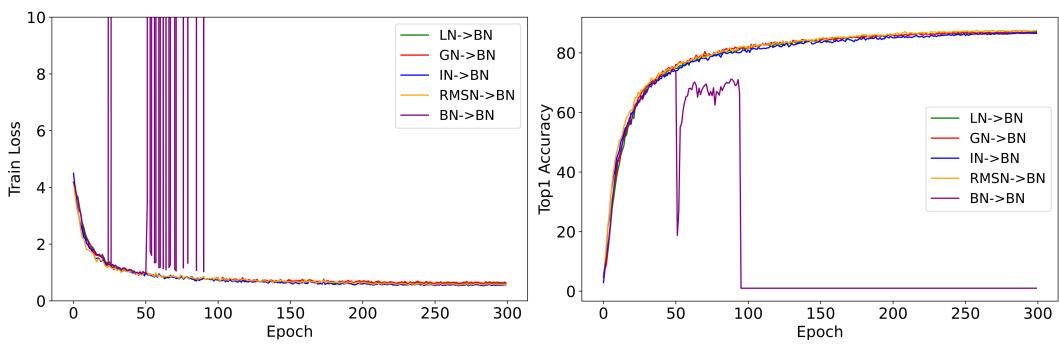


Figure 7: The results of combination other normalization with BN. From the decreasing loss and increasing accuracy, BN improves the performance when combined with other normalization layers, but reduces performance if used alone.

5 CONCLUSIONS AND FUTURE WORK

In this study, we conducted a comprehensive investigation into the types, positions, and combinations of normalization layers within the Mamba architecture. Our findings reveal that applying normalization after the S6 Modules enhances training stability by mitigating large variations in weight

486 norms. Moreover, specific combinations of normalization techniques not only stabilize the training
 487 process but also lead to significant improvements in model performance.
 488

489 We also proposed an intuition for selecting and combining normalization layers to optimize the training
 490 of Mamba and other deep architectures. This intuition provides valuable insights for managing
 491 weight domination and addressing the challenges of training large-scale neural networks. Future
 492 research will focus on extending this intuition to more complex models and tasks, to refine normal-
 493 ization strategies to enhance both efficiency and robustness further.
 494

495 The Mamba2 architecture has recently been introduced, and Dao & Gu (2024) found that training
 496 Mamba2 is less stable compared to Mamba1. Building upon our research into the normalization of
 497 Mamba Blocks, future work could further explore the Mamba2 architecture to address its stability
 498 challenges. Additionally, we hope our findings can be extended to other deep learning architectures,
 499 offering valuable guidance for designing more complex models.
 500

REFERENCES

- 501 Jimmy Lei Ba. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016.
 502
- 503 Tri Dao and Albert Gu. Transformers are ssms: Generalized models and efficient algorithms through
 504 structured state space duality. *ICML*, 2024.
- 505 Guillaume Desjardins, Karen Simonyan, Razvan Pascanu, et al. Natural neural networks. *Advances*
 506 *in neural information processing systems*, 28, 2015.
- 507 Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. *arXiv*
 508 *preprint arXiv:2312.00752*, 2023.
- 509 Albert Gu, Karan Goel, and Christopher Ré. Efficiently modeling long sequences with structured
 510 state spaces. *arXiv preprint arXiv:2111.00396*, 2021.
- 511 James D Hamilton. State-space models. *Handbook of econometrics*, 4:3039–3080, 1994.
- 512 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing
 513 human-level performance on imagenet classification. In *Proceedings of the IEEE international*
 514 *conference on computer vision*, pp. 1026–1034, 2015.
- 515 Geoffrey E Hinton and Ruslan R Salakhutdinov. Reducing the dimensionality of data with neural
 516 networks. *science*, 313(5786):504–507, 2006.
- 517 Lei Huang, Yi Zhou, Fan Zhu, Li Liu, and Ling Shao. Iterative normalization: Beyond standardiza-
 518 tion towards efficient whitening. In *Proceedings of the IEEE/CVF conference on computer vision*
 519 *and pattern recognition*, pp. 4874–4883, 2019.
- 520 Lei Huang, Jie Qin, Li Liu, Fan Zhu, and Ling Shao. Layer-wise conditioning analysis in exploring
 521 the learning dynamics of dnns. In *Computer Vision–ECCV 2020: 16th European Conference,*
 522 *Glasgow, UK, August 23–28, 2020, Proceedings, Part II 16*, pp. 384–401. Springer, 2020.
- 523 Lei Huang, Jie Qin, Yi Zhou, Fan Zhu, Li Liu, and Ling Shao. Normalization techniques in training
 524 dnns: Methodology, analysis and application. *IEEE transactions on pattern analysis and machine*
 525 *intelligence*, 45(8):10173–10196, 2023.
- 526 Xun Huang and Serge Belongie. Arbitrary style transfer in real-time with adaptive instance normal-
 527 ization. In *Proceedings of the IEEE international conference on computer vision*, pp. 1501–1510,
 528 2017.
- 529 Sergey Ioffe. Batch normalization: Accelerating deep network training by reducing internal covari-
 530 ate shift. *arXiv preprint arXiv:1502.03167*, 2015.
- 531 Diederik P Kingma. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*,
 532 2014.
- 533 Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convo-
 534 tional neural networks. *Advances in neural information processing systems*, 25, 2012.

540 Shen Li, Yanli Zhao, Rohan Varma, Omkar Salpekar, Pieter Noordhuis, Teng Li, Adam Paszke, Jeff
541 Smith, Brian Vaughan, Pritam Damania, et al. Pytorch distributed: Experiences on accelerating
542 data parallel training. *Proceedings of the VLDB Endowment*, 13(12), 2020.

543
544 Aobo Liang, Xingguo Jiang, Yan Sun, and Chang Lu. Bi-mamba4ts: Bidirectional mamba for time
545 series forecasting. *arXiv preprint arXiv:2404.15772*, 2024.

546 Yue Liu, Yunjie Tian, Yuzhong Zhao, Hongtian Yu, Lingxi Xie, Yaowei Wang, Qixiang Ye, and
547 Yunfan Liu. Vmamba: Visual state space model 2024. *CVPR*, 2024.

548
549 Ilya Loshchilov and Frank Hutter. Sgdr: Stochastic gradient descent with warm restarts. *Learning*,
550 10, 2016.

551 Xiong Luo, Xiaohui Chang, and Xiaojuan Ban. Regression and classification using extreme learning
552 machine based on l1-norm and l2-norm. *Neurocomputing*, 174:179–186, 2016.

553
554 Vardan Papyan. The full spectrum of deepnet Hessians at scale: Dynamics with SGD training and
555 sample size. *arXiv preprint arXiv:1811.07062*, 2018.

556 Yi Tay, Mostafa Dehghani, Dara Bahri, and Donald Metzler. Long range arena: A benchmark for ef-
557 ficient transformers. In *Proceedings of the International Conference on Learning Representations*
558 (*ICLR*), 2021.

559
560 Jiaxi Wang, Ji Wu, and Lei Huang. Understanding the failure of batch normalization for transformers
561 in nlp. *Advances in Neural Information Processing Systems*, 35:37617–37630, 2022.

562 Yuxin Wu and Kaiming He. Group normalization. In *Proceedings of the European conference on*
563 *computer vision (ECCV)*, pp. 3–19, 2018.

564 Biao Zhang and Rico Sennrich. Root mean square layer normalization. *Advances in Neural Infor-*
565 *mation Processing Systems*, 32, 2019.

566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593

594 **A TRAINING STABILITY**
 595

596 Given a normed vector space V , we refer to the L_1 and L_2 norm to be the special cases of the
 597 following general L_p norm of a give vector $x \in V$, by setting $p = 1$ and $p = 2$:
 598

599
$$\|x\|_{L^p} = \left(\sum_{j=1}^{\infty} |\xi_j|^p \right)^{1/p}$$

 600
 601
 602

603 We evaluate the training stability using L_2 norms. We analyzed L2 norm of the weight matrix of the
 604 whole Mamba Block,including in_projection layer,conv1d layer,S6 Module,out_projection layer.
 605

606 **B MAMBA ARCHITECTURE**
 607

608 State space modelling is a method of describing and analysing a system based on matrix theory.
 609 Introducing state variables can get more in-depth information about the system. SSM uses first-
 610 order differential equations to map the input function x_t to the output function y_t through hidden
 611 state h_t , defined as follows:
 612

613
$$h_t = \mathbf{A}h_{t-1} + \mathbf{B}x_t, \quad y_t = \mathbf{C}h_t + \mathbf{D}x_t \quad (2)$$

614 where state transition matrix $\mathbf{A} \in \mathbb{R}^{N \times N}$,input matrix $\mathbf{B} \in \mathbb{R}^{N \times C}$,output matrix $\mathbf{C} \in \mathbb{R}^{C \times N}$ and
 615 forward channel transition matrix $\mathbf{D} \in \mathbb{R}^{C \times C}$. The variables N and C refer to the hidden state
 616 and dimension factors, respectively. Continuous parameters \mathbf{A}, \mathbf{B} can be discretized by a first-order
 617 difference method or a bilinear transformation method as follows $\bar{\mathbf{A}}, \bar{\mathbf{B}}$, with the sampling interval
 618 Δ . Details of the process can be found in (Gu et al., 2021). We give the discretization result directly
 619 here. What should be noted is that the hidden state update mechanism of SSM is similar to Recurrent
 620 Neural Network (RNN), which receives current time step input x_t and the previous time step hidden
 621 state h_{t-1} and computes the current time step hidden state h_t .

622 The key design principle of Mamba lies in the introduction of a selective mechanism to parame-
 623 terize the transition matrices $\mathbf{A}, \mathbf{B}, \mathbf{C}, \Delta$ in a data-driven manner. In Mamba1, these matrices are
 624 defined as functions of the input embedding features, allowing them to adapt dynamically to the data
 625 context through a hardware-aware parallel computing algorithm. This enables efficient processing
 626 of long sequences with improved computational throughput. In contrast, Mamba2 leverages semi-
 627 differentiable matrix factorization to compute the hidden state space efficiently, ensuring that any
 628 state space model with state size N and sequence length L can be computed in time $O(TN)$, which
 629 means that Mamba2 preserves the controllability of the state space across various sequence lengths.
 630

631 **C IMPLEMENTATION DETAILS**
 632

633 **C.1 LISTO EXPERIMENT**

634 **Dataset** The Long Range Arena (LRA) benchmark (Tay et al., 2021) is designed to evaluate the
 635 ability of models to capture long-range dependencies across various tasks. One of the tasks included
 636 in the LRA benchmark is the *ListOps* task, which is particularly useful for assessing a model’s
 637 capability to handle hierarchical structures and process long-range sequences.
 638

639 In the *ListOps* task, the input sequence is represented as a nested list of operations. These symbolic
 640 expressions involve various operations, such as `max`, `min`, and `median`, applied to integers. The
 641 task requires the model to parse the input sequence, resolve the nested operations, and compute the
 642 final result. For example, a typical input could be:
 643

644
$$\max(2, \min(3, 9), \max(4, 5))$$

 645

646 In this example, the operations `max`, `min`, and `max` are applied in a hierarchical structure. The
 647 model’s objective is to correctly parse the operations and compute the final result, which in this case
 648 is 5.

The standard *ListOps* task involves input sequences of length 1000. However, to further test Mamba’s capacity to handle long-range dependencies, we used an enhanced dataset, *ListOps-New*, which includes input sequences of length 2000. This extended version allows us to assess better the limits of Mamba’s ability to manage and process extremely long sequences with complex, hierarchical operations.

Experiment Details All experiments were implemented in Pytorch(Listo (Tay et al., 2021)) and conducted on 12 NVIDIA 4090 24GB GPU using DDP (Li et al., 2020). We set the initial learning rate as 0.0001 and used the ADAMW optimizer (Kingma, 2014) for model optimization. The batch size was set to 32. To maintain stable training, we take the cosine_warmup scheduler (Loshchilov & Hutter, 2016). Detailed model configuration information is presented in Table 6

Table 6: Model Configuration Details

Encoder	Value
encoder_name	position
encoder_dropout	0.0
Layer	Value
layer_name	mamba
causal	false
Parameter	Value
dropout	0.0
n_layers	4
d_model	128
ss_state	64
d_conv	4
expand	2
epoch	30

Complete Result of Combanation We present the results of all possible combinations of commonly used normalization techniques in Table 7. We observed that combinations involving BN and GN often yield better performance. The combinations with GN tend to outperform those with BN, which we hypothesize is due to the inconsistencies between training and inference in BN, leading to performance degradation. The specific reasons behind this require further investigation in future studies.

Table 7: Accuracy Comparison for Various Normalization Combinations: The first arrow indicates a comparison with the combination in the first column, while the second arrow represents a comparison with the combination in the second column. Comparisons with close values are not marked.

Normalization	Accuracy	Normalization	Accuracy	Normalization	Accuracy	Normalization	Accuracy
BN → BN	0.439	IN → IN	0.406	BN → IN	0.469(↑)	IN → BN	0.376(↓)
BN → BN	0.439	RMSN → RMSN	0.401	BN → RMSN	0.401(↑)	RMSN → BN	0.448(↑)
BN → BN	0.439	LN → LN	0.381	BN → LN	0.386(↓)	LN → BN	0.418(↓)
BN → BN	0.439	GN → GN	0.414	BN → GN	0.401(↓)	GN → BN	0.429(↑)
IN → IN	0.406	LN → LN	0.381	IN → LN	0.398(↑)	LN → IN	0.412(↑)
IN → IN	0.406	RMSN → RMSN	0.401	IN → RMSN	0.392(↓)	RMSN → IN	0.372(↓)
IN → IN	0.406	GN → GN	0.414	IN → GN	0.468(↑)	GN → IN	0.448(↑)
LN → LN	0.381	RMSN → RMSN	0.401	LN → RMSN	0.393(↓)	RMSN → LN	0.400(↑)
LN → LN	0.381	GN → GN	0.414	LN → GN	0.404(↑)	GN → LN	0.425(↑)
RMSN → RMSN	0.401	GN → GN	0.414	RMSN → GN	0.423(↑)	GN → RMSN	0.416(↑)

C.2 IMAGENET EXPERIMENT

Experiment Details All experiments were implemented in Pytorch(Listo (Liu et al., 2024)) and conducted on 12 NVIDIA 4090 24GB GPU using DDP (Li et al., 2020). We set the initial learning rate as 0.001 and used the ADAMW optimizer (Kingma, 2014) for model optimization. The batch

size was set to 256. To maintain stable training, we take the cosine_warmup scheduler (Loshchilov & Hutter, 2016).

Result Showcase In Figure 8, we present the training loss and validation accuracy curves during the training of the VMamba model using different normalization layers and positioning them in different locations.

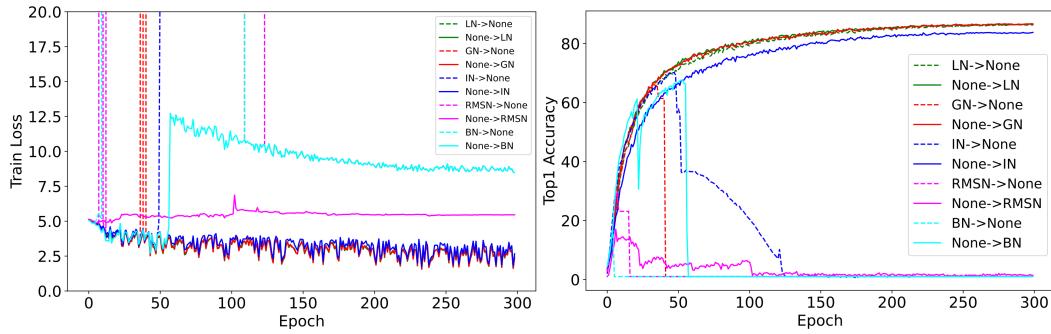


Figure 8: Impact of using five typical Normalization(BN, LN, IN, GN, RMS) in before or after S6 module. The solid line represents the results where normalization is applied after the S6 Module, while the dashed line represents the results where normalization is applied before the S6 Module.

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