

# MIX-LN: UNLEASHING THE POWER OF DEEPER LAYERS BY COMBINING PRE-LN AND POST-LN

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Large Language Models (LLMs) have achieved remarkable success, yet recent findings reveal that their deeper layers often contribute minimally and can be pruned without affecting overall performance. While some view this as an opportunity for model compression, we identify it as a training shortfall rooted in the widespread use of Pre-Layer Normalization (Pre-LN). We demonstrate that Pre-LN, commonly employed in models like GPT and LLaMA, leads to diminished gradient norms in its deeper layers, reducing their effectiveness. In contrast, Post-Layer Normalization (Post-LN) preserves larger gradient norms in deeper layers but suffers from vanishing gradients in earlier layers. To address this, we introduce MIX-LN, a novel normalization technique that combines the strengths of Pre-LN and Post-LN within the same model. MIX-LN applies Post-LN to the earlier layers and Pre-LN to the deeper layers, ensuring more uniform gradient norms across layers. This allows all parts of the network—both shallow and deep layers—to contribute effectively to training. Extensive experiments with various model sizes demonstrate that MIX-LN consistently outperforms both Pre-LN and Post-LN, promoting more balanced, healthier gradient norms throughout the network, and enhancing the overall quality of LLM pre-training. Furthermore, we demonstrate that models pre-trained with MIX-LN learn better compared to those using Pre-LN or Post-LN during supervised fine-tuning, highlighting the critical importance of high-quality deep layers. By effectively addressing the inefficiencies of deep layers in current LLMs, MIX-LN unlocks their potential, enhancing model capacity without increasing model size. Our code is submitted.

## 1 INTRODUCTION

Large Language Models (LLMs) have ushered in a new era of artificial intelligence by demonstrating unprecedented capabilities in understanding and generating human-like text (Brown, 2020; Achiam et al., 2023; Touvron et al., 2023; Dubey et al., 2024). Trained on vast datasets that span multiple languages and topics, LLMs are driving advancements across industries and academia, enhancing human-computer interactions, and fostering innovation in previously unimaginable ways.

Recent studies reveal a critical observation regarding the effectiveness of deeper layers in LLMs, particularly those beyond the middle layers. It has been shown that these deeper layers can often be pruned significantly (Yin et al., 2023), or even removed entirely (Gromov et al., 2024; Men et al., 2024), without notably affecting the model’s overall capabilities. Moreover, Li et al. (2024) demonstrated that deeper layers contribute minimally to performance during fine-tuning, further questioning their importance. Unfortunately, this finding has been largely overlooked by the research community, where many see it primarily as an opportunity for model compression (Siddiqui et al., 2024; Zhong et al., 2024; Sreenivas et al., 2024), rather than recognizing it as a potential shortfall in the training process.

In this paper, we seek to challenge the prevailing notion that deeper layers in LLMs are of lesser significance. The training of LLMs is extraordinarily resource-intensive, often requiring thousands of GPUs or TPUs and several months of computation on vast datasets. For example, the training of GPT-3 reportedly incurred millions of dollars in computational costs. The underutilization of deeper layers leads to inefficiencies, squandering resources that could otherwise be leveraged to enhance model performance. Ideally, all layers in a model should be well-trained, with sufficient diversity

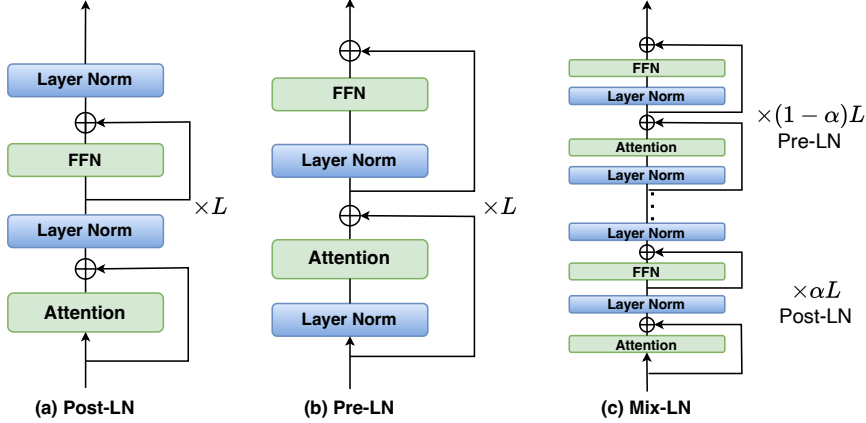


Figure 1: (a) Post-LN layer; (b) Pre-LN layer; (c) Mix-LN layer.

in features from layer to layer, to maximize the utility of the network’s parameters (Yang et al., 2023). This makes it crucial to investigate the root causes of this underutilization and to develop strategies that fully capitalize on the potential of deeper layers, ensuring that the overall architecture is optimized for both performance and efficiency.

We hypothesize that the inefficiency of deeper layers in LLMs primarily stems from the choice of Layer Normalization. Specifically, Pre-Layer Normalization (Pre-LN) (Dai, 2019; Baevski & Auli, 2018) tends to produce smaller gradients in deeper layers, thereby diminishing their effectiveness, while Post-Layer Normalization (Post-LN) (Ba, 2016) results in larger gradients in deeper layers but leads to gradient vanishing in earlier ones. Most state-of-the-art LLMs, like GPT, LLaMA, and Mistral, employ Pre-LN, which contributes to the widespread assumption that deeper layers are inherently less effective.

To validate this conjecture, we conduct experiments with the following two categories of LLMs and compare the effectiveness of layers across different depths in Pre-LN models and Post-LN models.

- **Open-weight large-scale LLMs:** We select LLaMa2-7B (Touvron et al., 2023) as a representative Pre-LN model and BERT-large (Devlin, 2018) as a Post-LN model to evaluate the quality of their layers. Our findings confirm that the deeper layers of LLaMa2-7B exhibit high similarity, with their removal leading to minimal impact compared to the early layers. In stark contrast, BERT shows higher similarity among its first half, which contributes less to the model’s output.
- **In-house small-scale LLMs:** To control for irrelevant confounding variables, we conduct a second set of experiments by training small-scale LLMs ourselves, ensuring that the only difference between the models is the choice of layer normalization. Consistent trends are observed in these experiments, reinforcing our earlier observations.

Building on these insights, we propose a novel normalization technique, dubbed `Mix-LN`, which synergizes Pre-LN and Post-LN to achieve more balanced and healthier gradient norms across the network. `Mix-LN` applies Post-LN to the earlier layers and Pre-LN to the deeper layers. The rationale behind this is that Post-LN enhances gradient flow in the deeper layers, while Pre-LN stabilizes gradients in the earlier layers. By employing Post-LN in the initial layers and Pre-LN in the later layers, `Mix-LN` promotes healthier gradient norms in the middle and deeper layers, fostering more balanced training across the entire network and ultimately improving the model’s overall performance.

Our extensive experiments, spanning models from 70M to 1B parameters, demonstrate that `Mix-LN` consistently outperforms Pre-LN, Post-LN, and their variants. `Mix-LN` not only avoids the training instability associated with Post-LN but also significantly improves the quality of deeper layers compared to Pre-LN, leading to better pre-training performance. Additionally, models pre-trained with `Mix-LN` demonstrate superior learning during supervised fine-tuning compared to those trained with Pre-LN or Post-LN, underscoring the importance of high-quality deep layers in LLMs.

## 2 HYPOTHESIS EVALUATION

In this section, we will evaluate our hypothesis that the inefficiency of deeper layers in LLMs stems from the choice of Pre-LN. The evaluation details are described as follows.

### 2.1 PRELIMINARIES: LAYER NORMALIZATION AND ITS GRADIENT

Figure 1 (a) and (b) illustrate Post-LN and Pre-LN Transformer architectures, respectively. Formally, let us define  $x$  as the input,  $\mathcal{F}(x)$  as either a FFN layer or a multi-head attention layer, and  $\text{LN}(\cdot)$  as the layer normalization. Post-LN applies  $\text{LN}(\cdot)$  after the residual addition:

$$\text{Post-LN}(x) = \text{LN}(x + \mathcal{F}(x)). \quad (1)$$

In contrast, Pre-LN applies  $\text{LN}(\cdot)$  before the residual addition:

$$\text{Pre-LN}(x) = x + \mathcal{F}(\text{LN}(x)). \quad (2)$$

It is well-known that shallow layers of Post-LN suffer from the gradient vanishing problem (Liu et al., 2020; Wang et al., 2024). Following Takase et al. (2022), we can calculate the derivatives of Equations (1) and (2), as follows:

$$\frac{\partial \text{Post-LN}(x)}{\partial x} = \frac{\partial \text{LN}(x + \mathcal{F}(x))}{\partial (x + \mathcal{F}(x))} \left( I + \frac{\partial \mathcal{F}(x)}{\partial x} \right), \quad (3)$$

$$\frac{\partial \text{Pre-LN}(x)}{\partial x} = I + \frac{\partial \mathcal{F}(\text{LN}(x))}{\partial \text{LN}(x)} \frac{\partial \text{LN}(x)}{\partial x}, \quad (4)$$

where  $I$  is the identity matrix. LN normalizes its input using the mean  $\mu$  and standard deviation  $\sigma$ ,  $\text{LN}(x) = \frac{x - \mu}{\sigma}$ . The derivative of  $\text{LN}(x)$  introduces a scaling factor  $\frac{1}{\sigma}$  term that accumulates as  $\prod_{l=1}^L \frac{1}{\sigma_l}$  over multiple layers  $L$ , which will reduce the gradient’s magnitude if  $\sigma > 1$ , which is often the case for Transformers<sup>1</sup>. Therefore, According to Eq. 3, such accumulation of gradient attenuation will cause gradient vanishing for the early layers of Post-LN. In contrast, Eq. 4 shows that the derivative of the residual connection is isolated from the term related to the derivative of LN, which prevents the vanishing gradient in early layers. On the other hand, Xiong et al. (2020) theoretically proved that the gradient norm of the LN depends on the norm of input as follows:

$$\left\| \frac{\partial \text{LN}(x)}{\partial x} \right\|_2 = \mathcal{O} \left( \frac{\sqrt{d}}{\|x\|_2} \right), \quad (5)$$

where  $d$  is the output dimension. Since Pre-LN does not normalize the residual connections, the output variance in Pre-LN models increases as the layers deepen, which in turn leads to diminished gradient norms in the deeper layers.

### 2.2 EVALUATION SETUP

**Methods:** Our evaluation methodology involves a comparative analysis of two models—one utilizing Pre-LN and the other employing Post-LN. By empirically assessing the effectiveness of layers across different depths in each model, we expect to see that Pre-LN models will exhibit a decrease in the effectiveness of deeper layers, whereas Post-LN models will show sustained or even improved quality in deeper layers.

**LLM Models:** To rigorously evaluate our hypothesis, we conduct experiments on two distinct categories of LLMs: (i) *Open-weight large-scale LLMs* and (ii) *In-house small-scale LLMs*. In the open-weight category, we select LLaMa2-7B (Touvron et al., 2023) as a representative Pre-LN model and BERT-large (Devlin, 2018) as a Post-LN model. However, these open-weight models differ not only in normalization but also in other factors such as training data, activation functions, and context length, complicating our ability to isolate the impact of normalization alone. To control for these confounding variables, we conduct a second set of experiments by training small-scale LLMs from scratch ourselves. The goal is to ensure that the only difference between the models is the choice of layer normalization. Specifically, we train LLaMa-130M models on the C4 dataset with either

<sup>1</sup>We also observe this in our experiments.

Pre-LN or Post-LN, using RMSNorm (Zhang & Sennrich, 2019) and SwiGLU activations (Shazeer, 2020), following Lialin et al. (2023b); Zhao et al. (2024). Please refer to Appendix A for more training configuration details.

**Evaluation Metrics:** A critical challenge in validating our hypothesis lies in defining and selecting robust metrics that capture the effectiveness of individual layers. In this study, we employ two metrics: (i) *Angular Distance* and (ii) *Performance Drop*, which provide a meaningful evaluation of the role and contribution of each layer. In addition, we report the *gradient norm* of each layer to demonstrate the effect of different layer normalization on the gradient flow.

(i) *Angular Distance*  $d(x^\ell, x^{\ell+n})$  is used in Gromov et al. (2024) to measure the angular distance between the input to layer  $\ell$  and the input to layer  $\ell + n$  on a neutral pre-training dataset. Formally, assuming  $x_T^\ell$  is the input to the layer  $\ell$ , and  $x_T^{\ell+n}$  is the input to the layer  $\ell + n$ , the angular distance between layers  $\ell$  and its subsequent  $n^{th}$  layer, i.e.,  $\ell + n$ , on a single token  $T$  is given by

$$d(x^\ell, x^{\ell+n}) = \frac{1}{\pi} \arccos \left( \frac{x_T^\ell \cdot x_T^{\ell+n}}{\|x_T^\ell\| \|x_T^{\ell+n}\|} \right) \quad (6)$$

where  $\|\cdot\|$  denotes the  $L^2$ -norm, and the factor of  $1/\pi$  scales  $d(x^\ell, x^{\ell+n})$  to the range  $[0, 1]$ . To eliminate the effect of randomness, the angular distance reported in this paper is averaged over 256K tokens from the C4 dataset. A smaller value of  $d(x^\ell, x^{\ell+n})$  indicates a shorter distance, meaning that the two vectors are more similar. Layers whose representations are extremely similar to their neighboring layers mean that they can be easily removed, and therefore their weights are less effective. Ideally, representation should change substantially from layer to layer in order to most effectively make use of the parameters of a network (Yang et al., 2023; Gromov et al., 2024).

(ii) *Performance Drop*  $\Delta P^{(\ell)}$  refers to the difference in the performance of an LLM before and after pruning the layer  $\ell$ . It quantifies the performance degradation caused by the removal of that layer. Formally, it can be defined as follows:

$$\Delta P^{(\ell)} = P_{\text{pruned}}^{(\ell)} - P_{\text{original}} \quad (7)$$

where  $P_{\text{original}}$  is the performance of the model without any pruning,  $P_{\text{pruned}}^{(\ell)}$  is the performance of the model after pruning layer  $\ell$ . A smaller value of  $\Delta P^{(\ell)}$  indicates that removing the layer causes minimal change to the model’s output, suggesting the layer is less important. Specifically, for LLaMA2-7B, we choose the commonly used MMLU (Hendrycks et al., 2020) as the evaluation task; for BERT-large, we opt for SQuAD v1.1 (Rajpurkar, 2016) as the evaluation task. Given the limited capacity of our in-house trained LLMs, we choose ARC-e (Clark et al., 2018) after supervised fine-tuning, instead of MMLU, for performance drop.

## 2.3 EVALUATION RESULTS

### 2.3.1 OPEN-WEIGHT LARGE-SCALE LLMs

Figure 2-(a, c) illustrate the metric values for BERT-Large. Both metrics indicate that, as a Post-LN model, the early layers of BERT-Large are less effective compared to the deeper layers. As shown in Figure 2-a, the first half of BERT-Large tends to have a smaller angular distance (more yellow) from neighboring layers than the second half. In particular, layers 3, 4, 9, 10, and 11 show a very high similarity to their subsequent layers. In Figure 2-c, the performance drop on SQuAD of removing an early layer is significantly smaller than the impact of removing a deeper layer. Intriguingly, removing layers 2 and 2 can even improve the performance.

In contrast, Figure 2-(b, d) display the metric values for LLaMa2-7B. As a Pre-LN model, the angular distance between neighboring layers decreases gradually (from purple to yellow) from the top layers to the 30th layer as illustrated in Figure 2-b. Notably, the deeper layers (20th to 30th) exhibit extremely small angular distances to their adjacent layers. This trend is consistent with the MMLU performance in Figure 2-d, where the removal of deeper layers results in almost negligible accuracy loss while removing early layers causes a substantial drop in accuracy.

In summary, we observe that the least effective layers in LLaMa2-7B are located in the deeper layers, whereas the early layers in BERT-Large are less effective than the deeper layers. The results from the category of open-weight large-scale LLMs strongly support our hypothesis, demonstrating a clear alignment with our expectations.

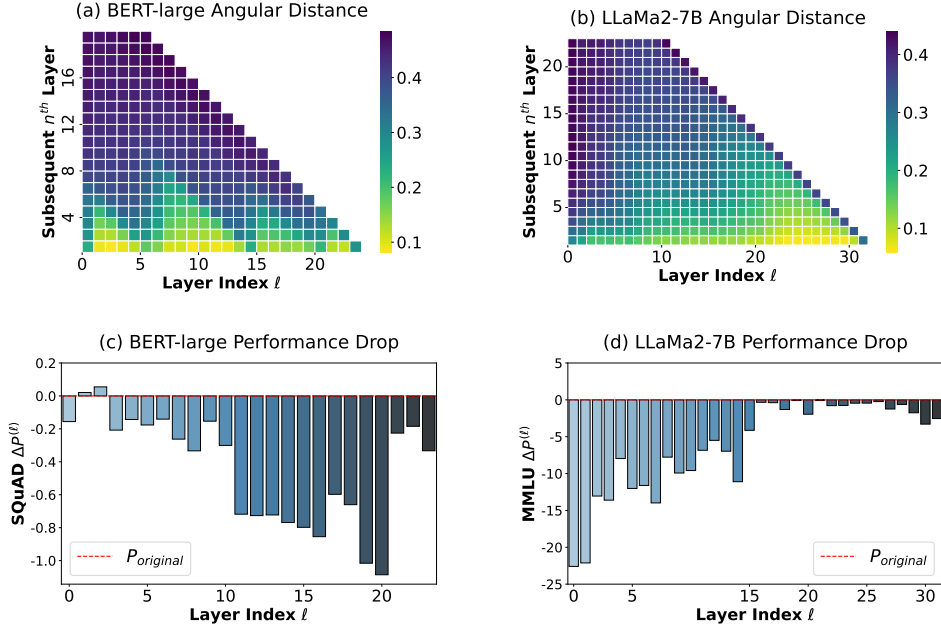


Figure 2: Results of open-weight large-scale LLMs. **Angular Distance (a, b):** Each column represents the angular distance from the initial layer  $\ell$  (x-axis) and its subsequent  $n^{th}$  layer (y-axis). The distance is scaled to the range  $[0, 1]$ , where yellow indicates smaller distances and purple indicates larger distances. **Performance Drop (c, d):** (c): SQuAD v1.1 performance drop of removing each single layer from BERT-large; (d): MMLU accuracy drop of removing each single layer from LLaMa2-7B.

### 2.3.2 IN-HOUSE SMALL-SCALE LLMs

Figure 3 illustrates all metric values for two LLaMa-130M models, where the only difference between them is the choice of layer normalization.

Figures 3-(a, b) show the Angular Distance for Post-LN and Pre-LN, respectively. Without the effects of other compounding factors, this comparison provides a clearer distinction between Post-LN and Pre-LN compared to open-weight large-scale LLMs. In Post-LN models, the most similar layers are concentrated in the early stages, with the first three layers showing particularly low distance. As the depth increases, the layers become increasingly distinctive. In contrast, the Pre-LN LLaMa-130M exhibits a gradual decrease in angular distance as depth increases, leading to highly similar deep layers. Figures 3-(d, e) further confirm this with the Performance Drop metric: removing early layers (e.g., 0-7 layers) in Post-LN results in minimal performance loss, while deeper layers (especially layers 9-11) are critical to preserving the original performance. However, Pre-LN LLaMa-130M exhibits the opposite trend, where removing most layers after the first layer causes negligible performance loss, indicating that they contribute little to the model’s output.

Figure 3-(c) shows the gradient norm of each layer for Post-LN and Pre-LN at the beginning of the training. The results perfectly align with our expectations: Post-LN leads to larger gradients in deeper layers but suffers from severe gradient vanishing in early layers, whereas Pre-LN maintains healthy gradient flow in early layers but diminishes in later layers.

With the consistent findings from both open-weight LLMs and our in-house LLMs, we can conclude that the widespread use of Pre-LN in LLMs is the root cause of the ineffectiveness of deep layers.

## 3 MIX-LAYER NORMALIZATION (MIX-LN)

Having validated our hypothesis that the use of Pre-LN is the root cause of the ineffectiveness of deeper layers, we propose **Mix-Layer Normalization (Mix-LN)**, a novel normalization strategy designed to enhance the effectiveness of both middle and deeper layers in LLMs.

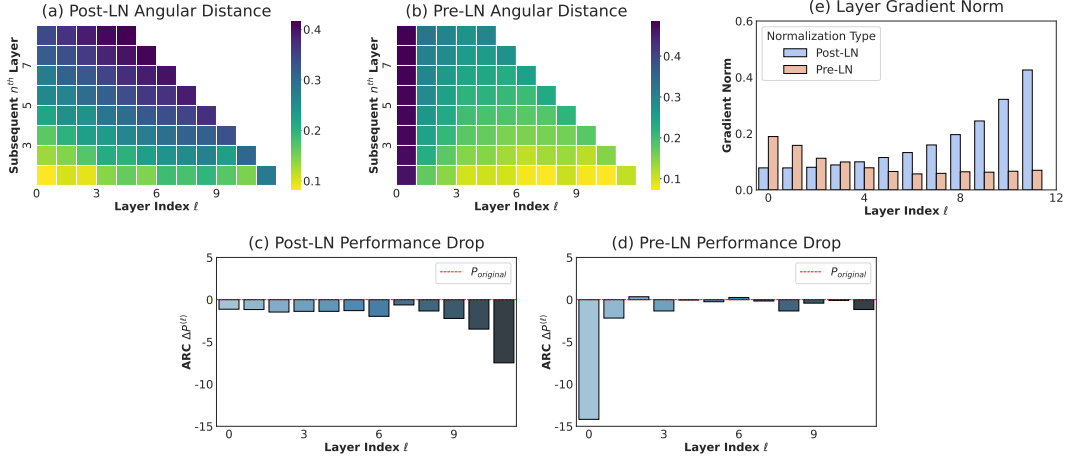


Figure 3: Results of in-house small-scale LLaMa-130M. **Angular Distance (a, b):** Each column represents the angular distance from the initial layer  $\ell$  (x-axis) and its subsequent  $n^{\text{th}}$  layer (y-axis). The distance is scaled to the range  $[0, 1]$ , where yellow indicates smaller distances and purple indicates larger distances. **Performance Drop (c, d):** ARC-e performance drop of removing each single layer from LLaMa-130M. **Gradient Norm (e):** Gradient norm of each layer in LLaMa-130M.

The key idea behind **Mix-LN** is to leverage the strengths of both Pre-LN and Post-LN. Post-LN has been shown to improve the effectiveness of deeper layers, while Pre-LN is more effective for earlier layers. Therefore, we propose to apply Post-LN to the initial layers and Pre-LN to the later layers, ensuring that the middle and deeper layers benefit from the advantages of both methods.

Formally, for an LLM with  $L$  layers, we apply Post-LN to the first  $\lfloor aL \rfloor$  layers and Pre-LN to the remaining  $\lceil (1 - a)L \rceil$  layers, where  $a \in [0, 1]$  is a hyperparameter controlling the transition point between the two normalization strategies. The functions  $\lfloor \cdot \rfloor$  and  $\lceil \cdot \rceil$  denote the floor and ceiling operations, respectively. Although the final layers may still experience smaller gradients due to the use of Pre-LN, the negative impact is substantially mitigated because the number of layers employing Pre-LN is reduced from  $L$  to  $\lceil (1 - a)L \rceil$ . This reduction improves gradient flow in the deeper layers, enhancing their effectiveness. Additionally, we expect that **Mix-LN** can alleviate training instability issues caused by Post-LN (Nguyen & Salazar, 2019; Wang et al., 2024), as reducing the number of layers using Post-LN leads to a smaller accumulation of gradient attenuation, according to the analysis in Section 2.1.

## 4 MAIN EXPERIMENTAL RESULTS

### 4.1 LLM PRE-TRAINING

In this section, we verify the effectiveness of **Mix-LN** by comparing it with various common normalization techniques, including Post-LN (Nguyen & Salazar, 2019), DeepNorm (Wang et al., 2024), and Pre-LN (Dai, 2019). Following Lialin et al. (2023a); Zhao et al. (2024), we conduct experiments using the LLaMA-based architecture with various sizes from 71M to 1B parameters, incorporating RMSNorm (Shazeer, 2020) and SwiGLU activations (Zhang & Sennrich, 2019). Models are trained with Adam Kingma (2014) using different learning rates based on model size: specifically, we use a learning rate of  $5e-3$  for models with 250M parameters and below, and a learning rate of  $5e-4$  for the 1B parameter model. All models of the same size are trained with identical configurations except for the normalization. To determine the optimal value for the hyperparameter  $\alpha$  in **Mix-LN**, we performed a small hyperparameter sweep using LLaMA-250M, as shown in Table 3. We found that  $\alpha = 0.25$  provided the best performance, and therefore, we applied this value across all model sizes.

Results are shown in Table 1. Post-LN generally yields the worst performance and even diverges with larger models, aligning with previous studies that indicate Post-LN suffers from training instability in Transformers (Xiong et al., 2020; Takase et al., 2022). DeepNorm, as a modified version of Post-LN, achieves comparable performance to Pre-LN with smaller model sizes; however, it also

Table 1: Perplexity comparison of various normalization methods across various LLaMA sizes.

Normalization	LLaMA-71M	LLaMA-130M	LLaMA-250M	LLaMA-1B
Training Tokens	1.1B	2.2B	3.9B	5B
Post-LN	35.18	32.18	1409.09	1411.54
DeepNorm	34.87	30.86	23.94	1410.94
Pre-LN	34.77	30.91	23.39	18.65
Mix-LN	<b>33.12</b>	<b>29.95</b>	<b>22.33</b>	<b>18.18</b>

experiences divergence during training with 1B parameter models. This observation confirms severe training instability of Post-LN, where gradients in early layers vanish, preventing proper model convergence. In contrast, Mix-LN consistently achieves the lowest perplexity across various model sizes. Mix-LN achieves a notable gain by 1.65 and 1.06 perplexity with LLaMA-71M and LLaMA-250M, respectively, compared to the widespread Pre-LN.

The above results clearly show that Mix-LN not only overcomes the instability of Post-LN but also enhances the model quality by combining the benefits of Pre-LN and Post-LN, making it an ideal choice for large-scale LLMs.

#### 4.2 SUPERVISED FINE-TUNING

Table 2: Fine-tuning performance of LLaMa with various normalizations.

Method	MMLU	BoolQ	ARC-e	PIQA	Hellaswag	OBQA	Winogrande	Avg.
<b>LLaMA-250M</b>								
Post-LN	22.95	37.83	26.94	52.72	26.17	11.60	49.56	32.54
DeepNorm	23.60	37.86	36.62	61.10	25.69	15.00	49.57	35.63
Pre-LN	24.93	38.35	40.15	63.55	26.34	16.20	49.01	36.93
Mix-LN	<b>26.53</b>	<b>56.12</b>	<b>41.68</b>	<b>66.34</b>	<b>30.16</b>	<b>18.00</b>	<b>50.56</b>	<b>41.34</b>
<b>LLaMA-1B</b>								
Post-LN	22.95	37.82	25.08	49.51	25.04	13.80	49.57	31.96
DeepNorm	23.35	37.83	27.06	52.94	26.19	11.80	49.49	32.67
Pre-LN	26.54	<b>62.20</b>	45.70	67.79	30.96	17.40	50.51	43.01
Mix-LN	<b>27.99</b>	61.93	<b>48.11</b>	<b>68.50</b>	<b>31.35</b>	<b>18.80</b>	<b>55.93</b>	<b>44.66</b>

We believe that the superior middle and deeper layers produced by Mix-LN are better equipped to learn during supervised fine-tuning. This advantage stems from the fact that these layers capture more diverse and rich features compared to those trained with Pre-LN. In complex downstream tasks, having access to a broad spectrum of features allows the model to make more nuanced predictions, leading to improved generalization.

To verify our conjecture, we follow Li et al. (2024) and fine-tune the models obtained in Section 4.1 on Commonsense170K (Hu et al., 2023), evaluating them on eight downstream tasks. As shown in Table 2, Mix-LN consistently outperforms other normalization techniques across all evaluated datasets. For the LLaMA-250M model, Mix-LN achieves a significant average gain of 4.26% and a 17.31% improvement on BoolQ compared to Pre-LN. Similar trends are observed with the larger LLaMA-1B model. Even though Mix-LN only slightly reduces perplexity by 0.25 compared to Pre-LN, it delivers substantial performance gains in supervised fine-tuning.

#### 4.3 SCALING UP TO 7B MODEL

Evaluating whether the benefits of Mix-LN scale to larger models, such as 7B parameters, is essential. To this end, we conducted experiments using the LLaMa-7B architecture, which features an embedding size of 4096 and 32 total layers, following the setup of Zhao et al. (2024). All training configurations were kept identical, with the exception of the layer normalization method. Due to computational constraints, we were able to complete only 13,000 steps of training. The training curve comparison is presented in Figure 4, where it is evident that Mix-LN consistently outperforms Pre-LN during the early training stages of LLaMa-7B.



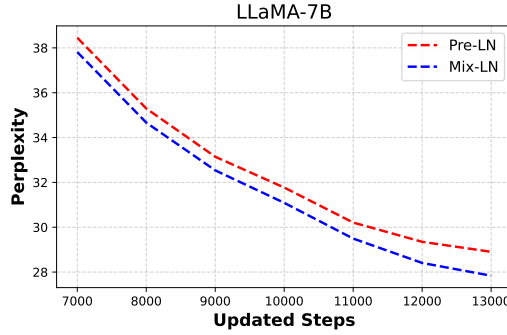
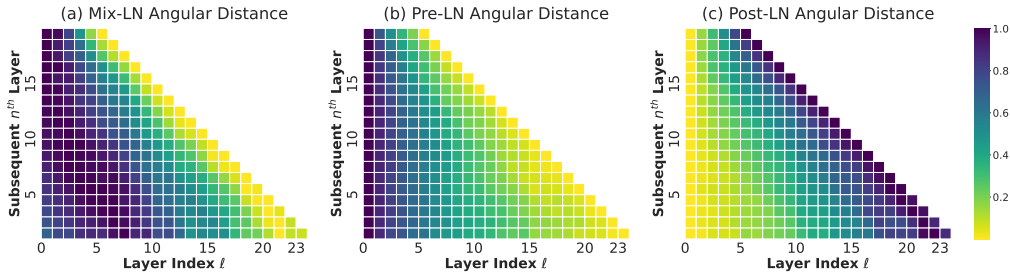


Figure 4: Training curve comparison between Mix-LN and Pre-LN with LLaMa-7B.

**Proper Post-LN ratio  $\alpha$ .** Mix-LN has a hyperparameter,  $\alpha$ , that controls the ratio of layers applying Post-LN. Specifically,  $\alpha = 0$  means Pre-LN is applied to all layers, while  $\alpha = 1$  corresponds to pure Post-LN. To determine the optimal Post-LN ratio, we conduct a sweep over the values [0, 12.5%, 25.0%, 33.0%, 41.7%, 50.0%, 75.0%, 100%] using LLaMA-250M on the C4 dataset. The results are shown in Table 3. As the normalization transitions from Pre-LN to Mix-LN, the model achieves progressively lower perplexity, reaching its best performance at  $\alpha = 0.25$ . Beyond this point, performance begins to decline, although it still surpasses that of pure Pre-LN until most layers apply Post-LN, where performance degrades significantly. Based on these results, we choose  $\alpha = 0.25$  for all model sizes, although we believe there is potential to further improve the performance of Mix-LN by searching for the optimal  $\alpha$  for each individual model.

Table 3: Perplexity of LLaMA-250M with various Post-LN ratios  $\alpha$ .

	Pre-LN	Mix-LN						Post-LN
Post-LN ratios $\alpha$	0	12.5%	25.0%	33.0%	41.7%	50.0%	75.0%	100%
Perplexity	23.39	22.37	<b>22.33</b>	22.83	22.80	22.81	23.64	32.18

Figure 5: Normalized angular distance from initial layer  $\ell$  (x-axis) with block size  $n$  (y-axis).

**Mix-LN promotes representation diversity across layers.** As we have claimed, our hybrid approach promotes a more balanced gradient flow throughout the entire network. To validate this, we report the angular distance of LLaMA-250M for Pre-LN, Post-LN, and Mix-LN in Figure 5. Following Gromov et al. (2024), we normalize each row to display the row-normalized angular distance between layer  $\ell$  (x-axis) and  $\ell + n$  (y-axis) for all possible  $\ell$ . Given block size  $n$ , the layers with the smallest distances are highlighted in the lightest yellow in each row. Notably, Mix-LN consistently exhibits larger distances (darker color) across layers compared to Pre-LN, except for the final two layers. This indicates that Mix-LN produces more diverse representations between layers than Pre-LN. In contrast, the smallest distances in Post-LN are concentrated in the early layers, reinforcing the notion that Post-LN tends to restrict representation diversity in deeper layers.

**Mix-LN enhances healthier gradient norms across all layers.** We compare the gradient norm of different LN at initialization in Figure 6. It demonstrates that Mix-LN maintains more consistent gradient norms across all layers. This balance results in a more uniform distribution of gradient



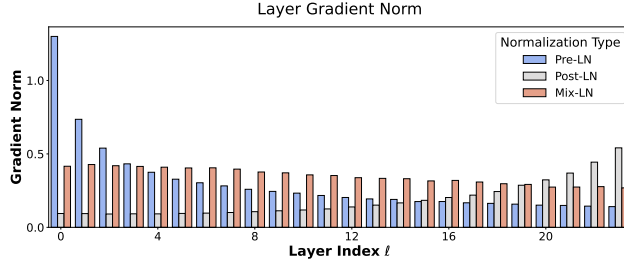


Figure 6: Layer gradient norm of LLaMA-250M with various normalization techniques.

norms across layers, allowing all parts of the network—both shallow and deep layers—to contribute effectively to model training.

## 5 RELATED WORK

### 5.1 NORMALIZATION IN LANGUAGE MODELS

Layer Normalization (LN), first proposed by Ba (2016), has become the de facto standard for normalizing activations in modern language models. It directly estimates normalization statistics from the summed inputs to neurons within a hidden layer, ensuring that the input distribution to each layer remains stable throughout training. In the original Transformer (Vaswani, 2017), LN was initially applied after the residual connection, a configuration known as Post-LN. However, subsequent studies (Baeviski & Auli, 2018; Dai, 2019; Nguyen & Salazar, 2019) found that placing LayerNorm before the residual connection (Pre-LN) results in more stable performance, especially in large language models (Brown, 2020; Touvron et al., 2023; Jiang et al., 2023). Xiong et al. (2020) theoretically demonstrated that Post-LN results in larger gradients near the output layer, making the use of warm-up essential to avoid instability is necessary. Conversely, Pre-LN scales down gradients with the depth of the model, which ensures more stable gradients during initialization. Our work builds upon Xiong et al. (2020), highlighting that while Pre-LN prevents instability by reducing gradient magnitudes, smaller gradients in deeper layers can diminish the effectiveness of the corresponding weights.

To improve the effectiveness of deeper layers in language models, various LN variants have been proposed. For instance, Wang et al. (2019) verified empirically that Post-LN suffers from gradient vanishing in deep Transformers, while Pre-LN facilitates stacking more layers. They consequently introduced dynamic linear combination of layers (DLCL), which connects all previous layers to improve trainability. Similar techniques have been employed in other works (Bapna et al., 2018; Dou et al., 2018). Liu et al. (2020) revealed that Post-LN has strong dependencies on the residual branch, often leading to instability. To address this, Adaptive Model Initialization (Admin) was introduced, which uses additional parameters to control residual dependencies in Post-LN, stabilizing training. DeepNorm (Wang et al., 2024) further improved the trainability of deep Transformers by upscaling the residual connection before applying LN, reducing model updates, and enabling deeper architectures. Additionally, Ding et al. (2021) proposed Sandwich LayerNorm, normalizing both the input and output of each transformer sub-layer. Takase et al. (2022) identified that Post-LN tends to preserve larger gradient norms in deeper layers, potentially leading to more effective training. To address the issue of gradient vanishing in early layers, they introduced B2T, a method that uses a residual connection to bypass all LN except the final one in each layer. We got inspiration from Takase et al. (2022), addressing the limitations of both Pre-LN and Post-LN by combining them. We study Scaled Initialization and Scaled Embed in Appendix B.

### 5.2 INEFFECTICACY OF DEEP LAYERS IN LLMs

The Inefficacy of deep layers in LLMs serves as a valid indicator for LLM pruning. Yin et al. (2023) demonstrated that the deeper layers of prominent LLMs like LLaMA and Mistral can be pruned more aggressively than earlier layers, without causing a significant drop in performance. Similarly, Gromov et al. (2024) and Men et al. (2024) further explored layer pruning, identifying the deeper layers of LLMs as typically less essential. Lad et al. (2024) observed that in models like Pythia and

GPT-2, deeper layers exhibit strong resilience to interventions, such as layer deletion or swapping. Our work shares similarities with Gromov et al. (2024) in applying angular distance to assess the effectiveness of layers. However, while they identify the inefficacy of deeper layers, they do not offer an explanation for this phenomenon nor propose a solution to address it.

While previous studies often view these characteristics of deeper layers as an opportunity for model compression (Siddiqui et al., 2024; Zhong et al., 2024; Sreenivas et al., 2024), we argue that this behavior reveals a deeper training shortfall—primarily due to the widespread use of Pre-LN. In response, we introduce `Mix-LN`, a novel method that enhances the effectiveness of deeper layers, ensuring that the entire architecture is more effectively trained and fully leverages the network’s parameters.

## 6 CONCLUSION

In this paper, we have addressed the inefficiencies of deep layers in LLMs by identifying the widespread use of Pre-LN as the root cause. Pre-LN leads to diminished gradients in deeper layers, reducing their effectiveness. While Post-LN preserves deeper gradients, it suffers from vanishing gradients in earlier layers. To resolve this, we introduced `Mix-LN`, a hybrid normalization technique that combines the strengths of both Pre-LN and Post-LN. By applying Post-LN to early layers and Pre-LN to deeper layers, `Mix-LN` achieves balanced gradient norms throughout the network, enabling more effective training. Our experiments show that `Mix-LN` consistently outperforms both Pre-LN and Post-LN, enhancing pre-training and fine-tuning performance without increasing model size. By fully utilizing the potential of deep layers, `Mix-LN` improves the overall capacity and efficiency of LLMs.

## REFERENCES

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- Jimmy Lei Ba. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016.
- Alexei Baevski and Michael Auli. Adaptive input representations for neural language modeling. *arXiv preprint arXiv:1809.10853*, 2018.
- Ankur Bapna, Mia Xu Chen, Orhan Firat, Yuan Cao, and Yonghui Wu. Training deeper neural machine translation models with transparent attention. *arXiv preprint arXiv:1808.07561*, 2018.
- Tom B Brown. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*, 2020.
- Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge. *arXiv preprint arXiv:1803.05457*, 2018.
- Zihang Dai. Transformer-xl: Attentive language models beyond a fixed-length context. *arXiv preprint arXiv:1901.02860*, 2019.
- Jacob Devlin. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- Ming Ding, Zhuoyi Yang, Wenyi Hong, Wendi Zheng, Chang Zhou, Da Yin, Junyang Lin, Xu Zou, Zhou Shao, Hongxia Yang, et al. Cogview: Mastering text-to-image generation via transformers. *Advances in neural information processing systems*, 34:19822–19835, 2021.
- Zi-Yi Dou, Zhaopeng Tu, Xing Wang, Shuming Shi, and Tong Zhang. Exploiting deep representations for neural machine translation. *arXiv preprint arXiv:1810.10181*, 2018.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.

- Andrey Gromov, Kushal Tirumala, Hassan Shapourian, Paolo Glorioso, and Daniel A Roberts. The unreasonable ineffectiveness of the deeper layers. *arXiv preprint arXiv:2403.17887*, 2024.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding. *arXiv preprint arXiv:2009.03300*, 2020.
- Zhiqiang Hu, Lei Wang, Yihuai Lan, Wanyu Xu, Ee-Peng Lim, Lidong Bing, Xing Xu, Soujanya Poria, and Roy Ka-Wei Lee. Llm-adapters: An adapter family for parameter-efficient fine-tuning of large language models. *arXiv preprint arXiv:2304.01933*, 2023.
- Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. Mistral 7b. *arXiv preprint arXiv:2310.06825*, 2023.
- Diederik P Kingma. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Vedang Lad, Wes Gurnee, and Max Tegmark. The remarkable robustness of llms: Stages of inference? *arXiv preprint arXiv:2406.19384*, 2024.
- Pengxiang Li, Lu Yin, Xiaowei Gao, and Shiwei Liu. Owlcore: Outlier-weighted layerwise sampled low-rank projection for memory-efficient llm fine-tuning. *arXiv preprint arXiv:2405.18380*, 2024.
- Vladislav Lialin, Sherin Muckatira, Namrata Shivagunde, and Anna Rumshisky. Relora: High-rank training through low-rank updates. In *The Twelfth International Conference on Learning Representations*, 2023a.
- Vladislav Lialin, Namrata Shivagunde, Sherin Muckatira, and Anna Rumshisky. Stack more layers differently: High-rank training through low-rank updates. *arXiv preprint arXiv:2307.05695*, 2023b.
- Liyuan Liu, Xiaodong Liu, Jianfeng Gao, Weizhu Chen, and Jiawei Han. Understanding the difficulty of training transformers. *arXiv preprint arXiv:2004.08249*, 2020.
- Xin Men, Mingyu Xu, Qingyu Zhang, Bingning Wang, Hongyu Lin, Yaojie Lu, Xianpei Han, and Weipeng Chen. Shortgpt: Layers in large language models are more redundant than you expect. *arXiv preprint arXiv:2403.03853*, 2024.
- Toan Q Nguyen and Julian Salazar. Transformers without tears: Improving the normalization of self-attention. *arXiv preprint arXiv:1910.05895*, 2019.
- P Rajpurkar. Squad: 100,000+ questions for machine comprehension of text. *arXiv preprint arXiv:1606.05250*, 2016.
- Tevan Le Scao, Thomas Wang, Daniel Hesslow, Lucile Saulnier, Stas Bekman, M Saiful Bari, Stella Biderman, Hady Elsahar, Niklas Muennighoff, Jason Phang, et al. What language model to train if you have one million gpu hours? *arXiv preprint arXiv:2210.15424*, 2022.
- Noam Shazeer. Glu variants improve transformer. *arXiv preprint arXiv:2002.05202*, 2020.
- Shoaib Ahmed Siddiqui, Xin Dong, Greg Heinrich, Thomas Breuel, Jan Kautz, David Krueger, and Pavlo Molchanov. A deeper look at depth pruning of llms. *arXiv preprint arXiv:2407.16286*, 2024.
- Sharath Turuvekere Sreenivas, Saurav Muralidharan, Raviraj Joshi, Marcin Chochowski, Mostafa Patwary, Mohammad Shoeybi, Bryan Catanzaro, Jan Kautz, and Pavlo Molchanov. Llm pruning and distillation in practice: The minitron approach. *arXiv preprint arXiv:2408.11796*, 2024.
- Sho Takase, Shun Kiyono, Sosuke Kobayashi, and Jun Suzuki. B2t connection: Serving stability and performance in deep transformers. *arXiv preprint arXiv:2206.00330*, 2022.
- Sho Takase, Shun Kiyono, Sosuke Kobayashi, and Jun Suzuki. Spike no more: Stabilizing the pre-training of large language models. *arXiv preprint arXiv:2312.16903*, 2023.

Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.

A Vaswani. Attention is all you need. *Advances in Neural Information Processing Systems*, 2017.

Hongyu Wang, Shuming Ma, Li Dong, Shaohan Huang, Dongdong Zhang, and Furu Wei. Deepnet: Scaling transformers to 1,000 layers. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024.

Qiang Wang, Bei Li, Tong Xiao, Jingbo Zhu, Changliang Li, Derek F Wong, and Lidia S Chao. Learning deep transformer models for machine translation. *arXiv preprint arXiv:1906.01787*, 2019.

Ruibin Xiong, Yunchang Yang, Di He, Kai Zheng, Shuxin Zheng, Chen Xing, Huishuai Zhang, Yanyan Lan, Liwei Wang, and Tieyan Liu. On layer normalization in the transformer architecture. In *International Conference on Machine Learning*, pp. 10524–10533. PMLR, 2020.

Greg Yang, Dingli Yu, Chen Zhu, and Soufiane Hayou. Tensor programs vi: Feature learning in infinite-depth neural networks. *arXiv preprint arXiv:2310.02244*, 2023.

Lu Yin, You Wu, Zhenyu Zhang, Cheng-Yu Hsieh, Yaqing Wang, Yiling Jia, Mykola Pechenizkiy, Yi Liang, Zhangyang Wang, and Shiwei Liu. Outlier weighed layerwise sparsity (owl): A missing secret sauce for pruning llms to high sparsity. *arXiv preprint arXiv:2310.05175*, 2023.

Biao Zhang and Rico Sennrich. Root mean square layer normalization. *Advances in Neural Information Processing Systems*, 32, 2019.

Jiawei Zhao, Zhenyu Zhang, Beidi Chen, Zhangyang Wang, Anima Anandkumar, and Yuandong Tian. Galore: Memory-efficient llm training by gradient low-rank projection. *arXiv preprint arXiv:2403.03507*, 2024.

Longguang Zhong, Fanqi Wan, Ruijun Chen, Xiaojun Quan, and Liangzhi Li. Blockpruner: Fine-grained pruning for large language models. *arXiv preprint arXiv:2406.10594*, 2024.

## A DETAILS OF EXPERIMENTS

### A.1 ARCHITECTURE AND HYPERPARAMETERS

We introduce details of the LLaMA architecture and hyperparameters used for pre-training following (Lialin et al., 2023a; Zhao et al., 2024). Table 4 shows the most hyperparameters of LLaMA models across model sizes. We use a max sequence length of 256 for all models, with a batch size of 512, with a batch size of 131K tokens. Learning rate warmup is applied to the first 10% of the training steps. We train models using Adam with a cosine annealing for the learning rate schedule, decaying to 10% of the initial learning rate. We use a learning rate of 5e-3 for models with 250M parameters and below, and a learning rate of 5e-4 for the 1B parameter model.

Table 4: Hyperparameters of LLaMA models used in this paper.

Params	Hidden	Intermediate	Heads	Layers	Steps	Data amount
71M	512	1368	8	12	10K	1.1 B
130M	768	2048	12	12	20K	2.2 B
250M	1024	2560	16	24	60K	3.9 B
1 B	2048	5461	24	32	100K	5.0 B

## B COMPATIBILITY TO ADVANCED

In this section, we also evaluate if `Mix-LN` can integrate well with the advanced techniques proposed to stabilize training. Specifically, we evaluate the commonly used Scaled Initialization (Nguyen & Salazar, 2019; Scao et al., 2022) that initializes  $W_2$  and  $W_0$  with a smaller normal distribution  $\mathcal{N}(0, \sqrt{2/5d}/\sqrt{2N})$  to stabilize training dynamics; and Scaled Embed (Takase et al., 2023) scales up embeddings to stabilize LayerNorm gradients. We observe that both Pre-LN and `Mix-LN` work effectively with Scaled Initialization. However, incorporating Scaled Embed on top of this setup leads to a degradation in performance.

Table 5: Perplexity of LLaMA-130M with various normalization methods with Scaled Initialization and Scaled Embed.

Normalization	Scaled Initialization	Scaled Embed	Perplexity
Pre-LN			32.18
Mix-LN			29.95
Pre-LN	✓		30.63
Mix-LN	✓		29.77
Pre-LN	✓	✓	31.28
Mix-LN	✓	✓	31.19