# Understanding the Difficulty of Training Transformers
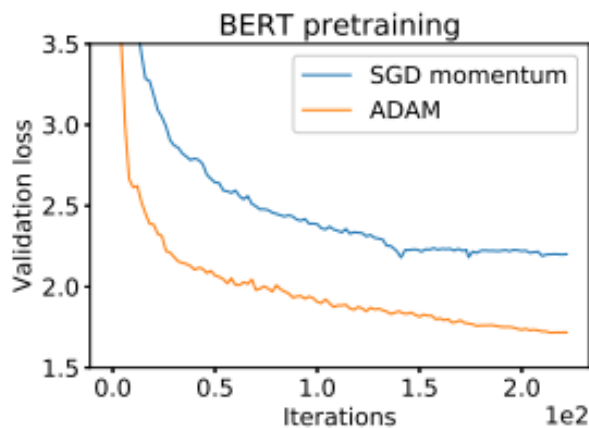
Liyuan Liu[†‡]  Xiaodong Liu[‡]  Jianfeng Gao[‡] Weizhu Chen[§] Jiawei Han[†]
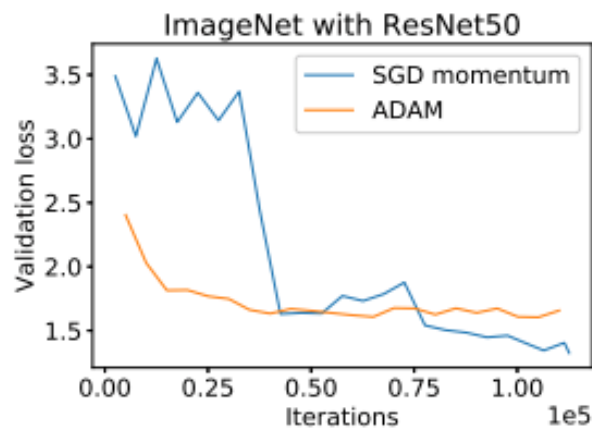
† University of Illinois at Urbana-Champaign
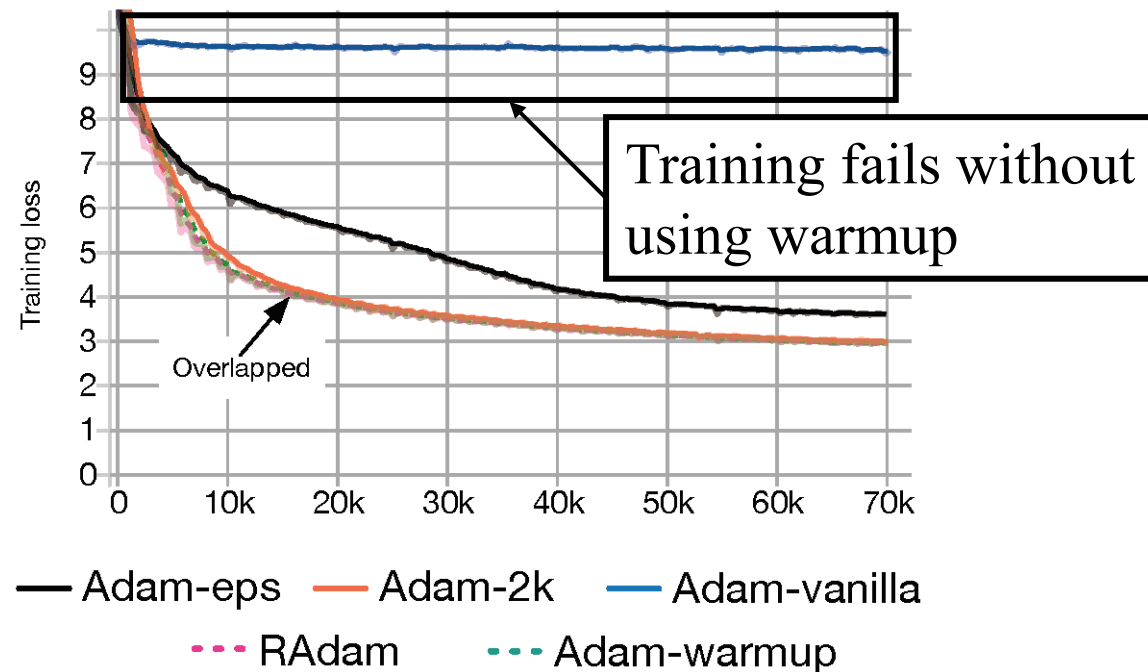‡ Microsoft Research
§ Microsoft Dynamics 365 AI

BERT pretraining

ImageNet with ResNet50

Training fails without using warmup

Overlapped

Adam-eps — Adam-2k — Adam-vanilla

RAdam — Adam-warmup

Although SGD is the canonical algorithm for conventional NNs, it fails to train Transformer effectively.

Removing the warmup phrase results in more serious consequences.

# Transformer requires non-trivial efforts

## What Complicates Transformer Training?

# Gradients Vanishing

Unbalanced gradients can hamper the training from the beginning and has been long regarded as the major reason destabilize model training.

Recent study shows that, even after introducing residual connections, the Transformer network still suffers from gradient vanishing.

Surprisingly, we find gradient vanishing is not the direct reason

Fixing the gradient vanishing issue alone cannot stabilize training.

Unbalanced gradients are largely handled by adaptive optimizers.

*Fixing Initialization Stabilizes the Transformer Training*

18-Layer Transformer
6-Layer Transformer

Pre-LN
Post-LN
Admin (Post-LN)

Without Admin, 6-layer Post-LN converges, and 18-layer diverges.

Dev PPL on WMT'14 En-De

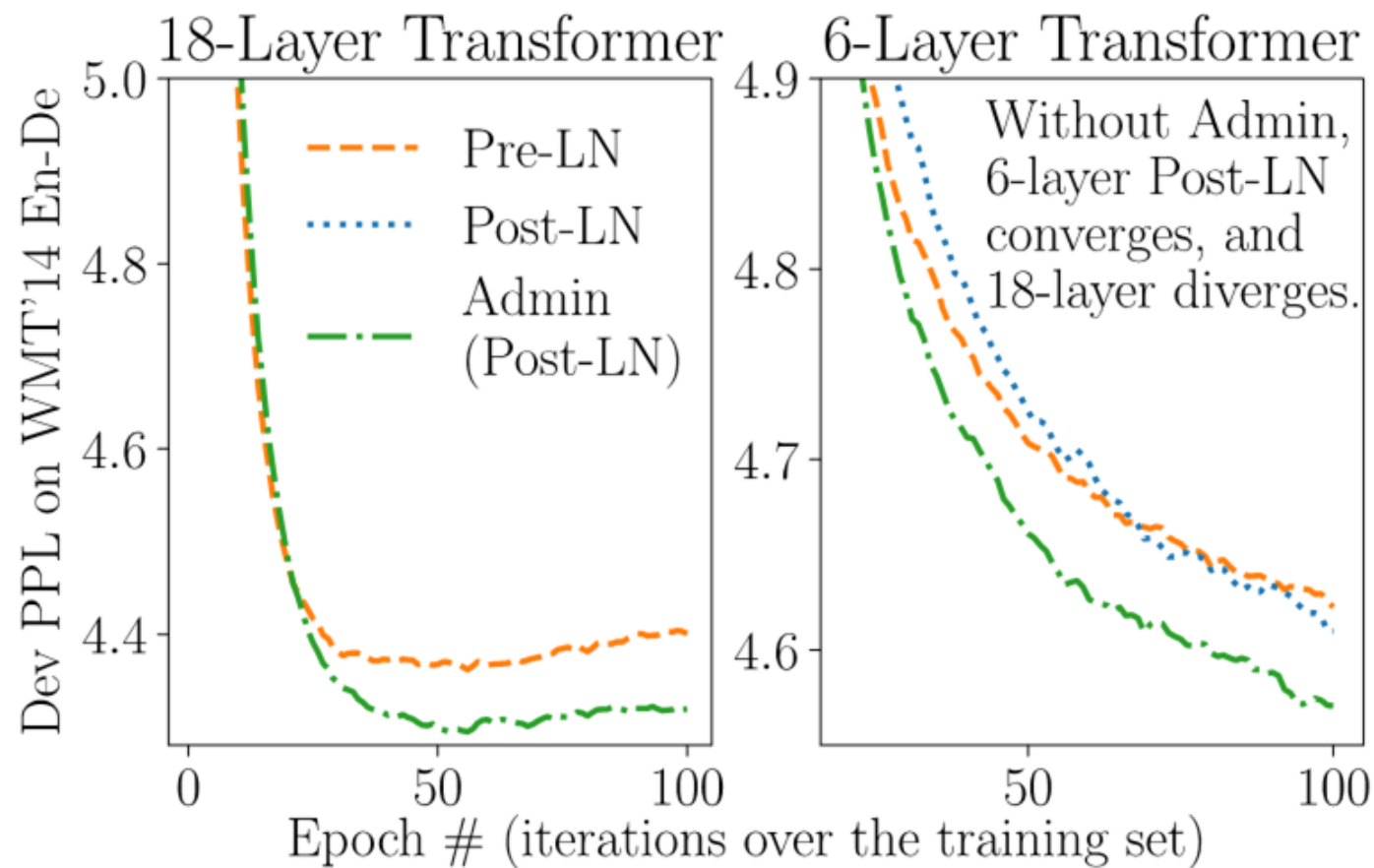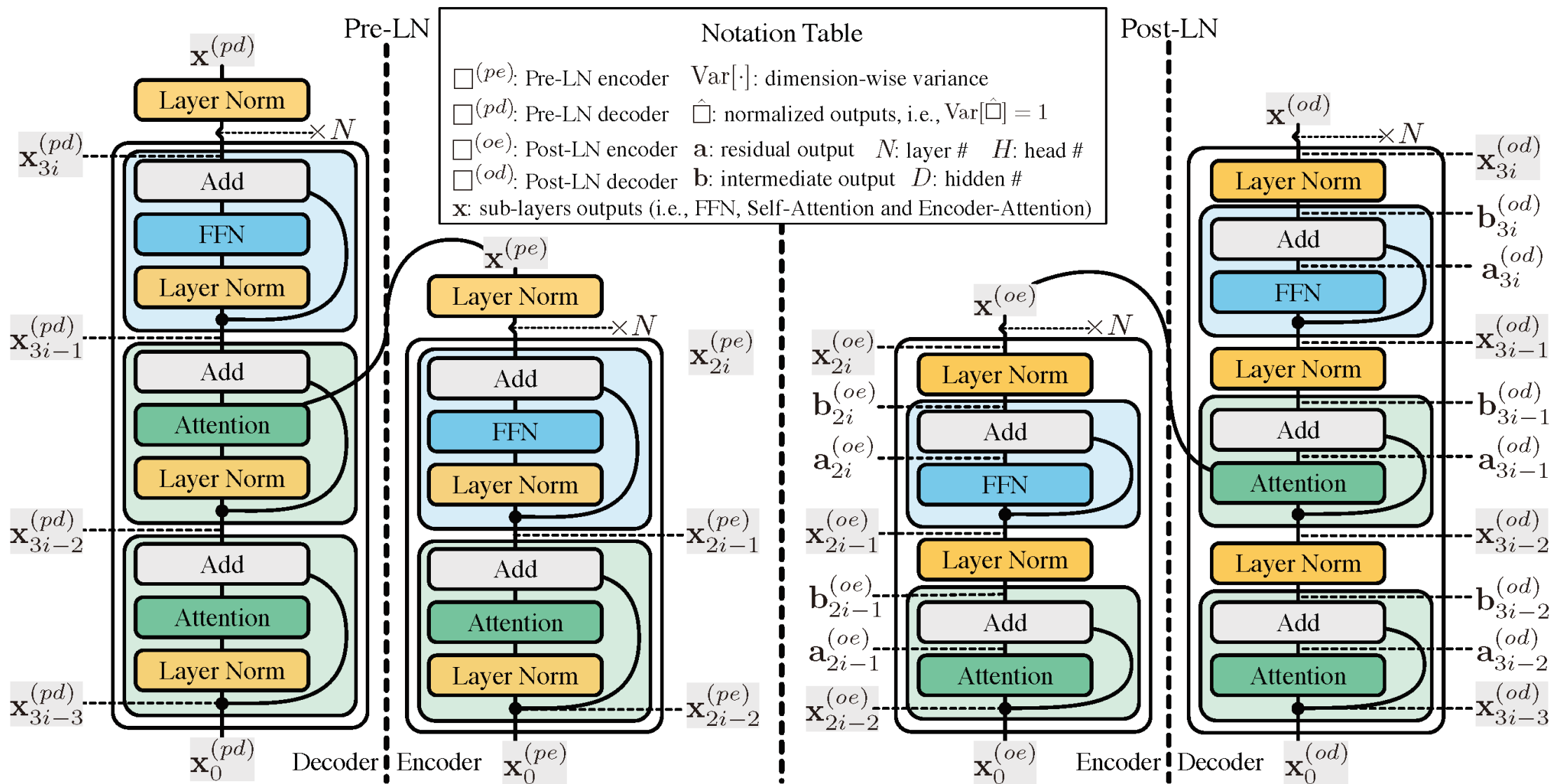Epoch # (iterations over the training set)

# Table of Contents

1. Gradient Vanishing and Transformer Training

2. Amplification Effect for Transformer Training
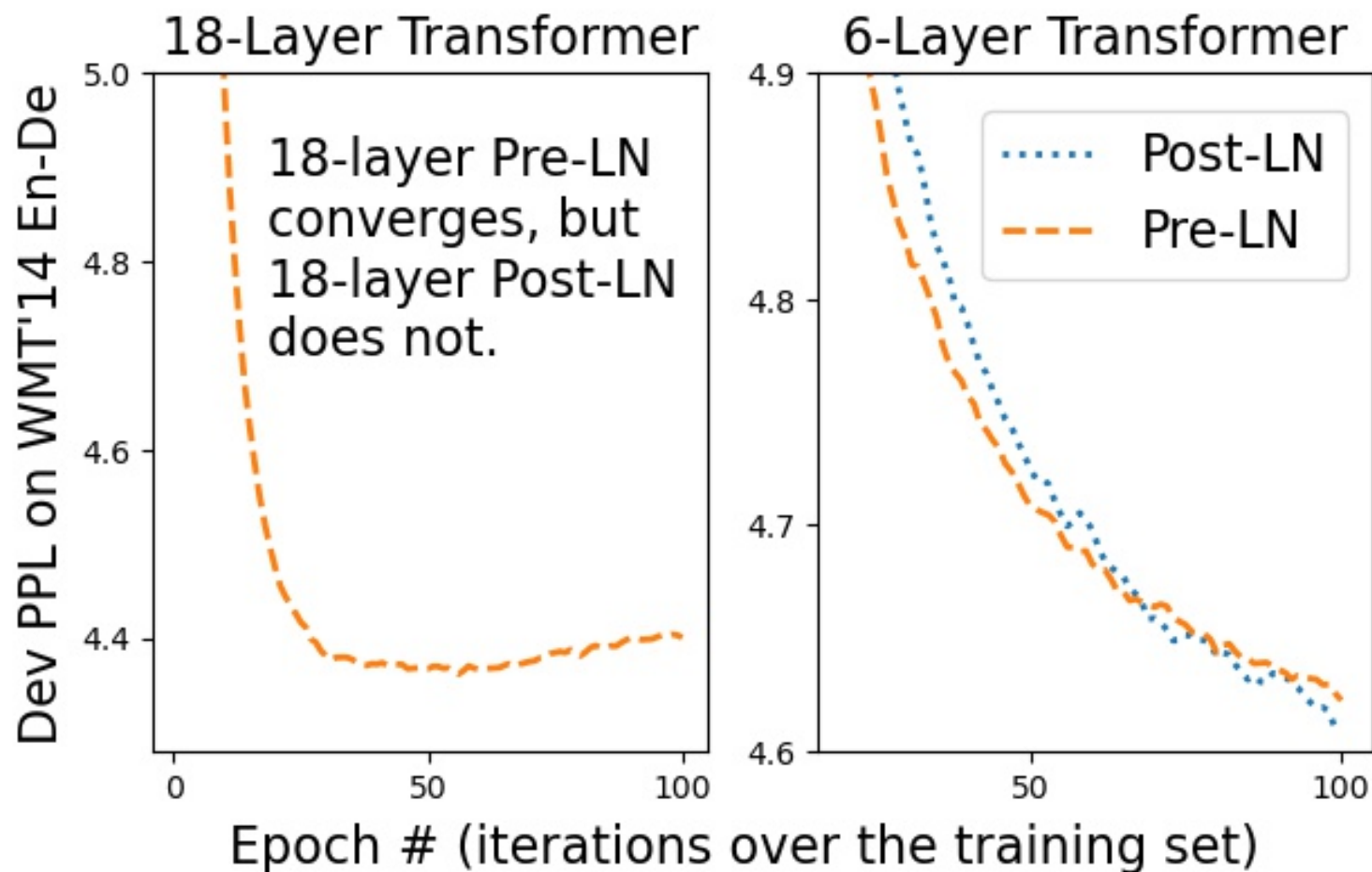
3. ADMIN (Adaptive Model Initialization)

Difference between the Pre-LN and the Post-LN:

Pre-LN variants are more robust.

Post-LN variants have a larger potential.
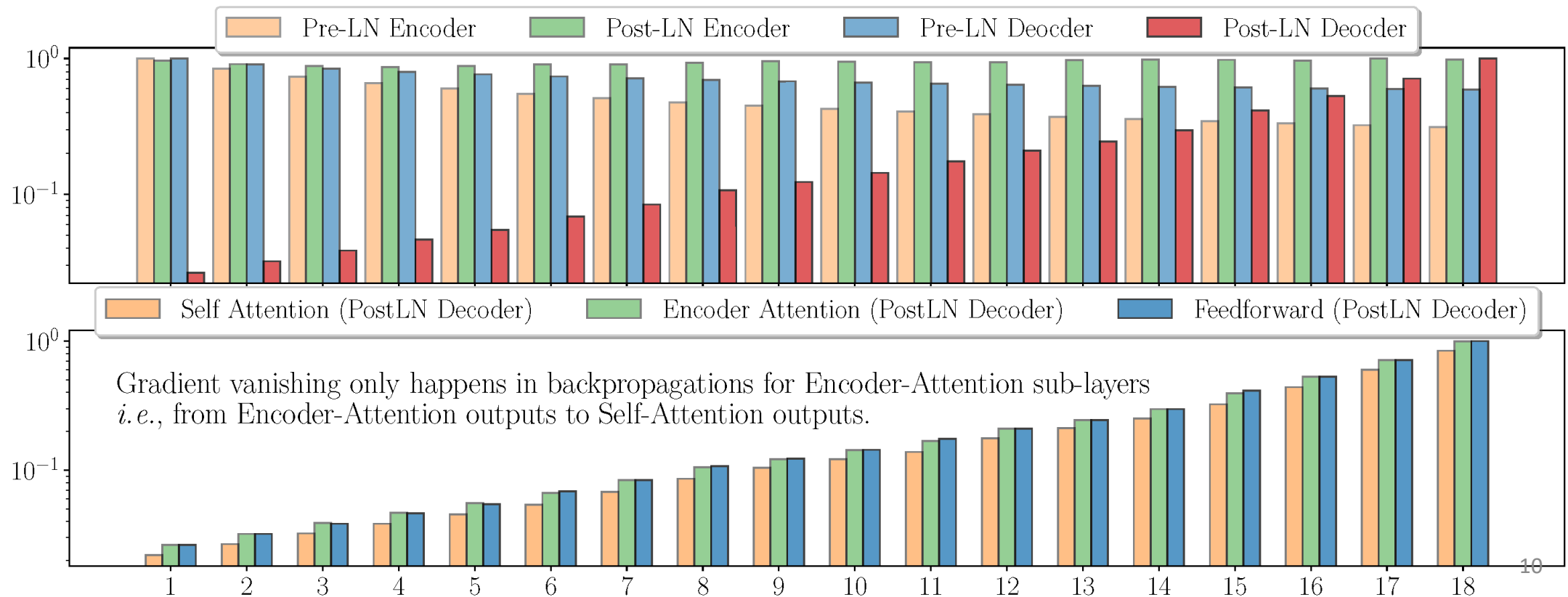
Our study
starts from ...

# Gradient Vanishing in Transformer

Since Post-LN suffers from gradient vanishing and is not stable, it is long believed that the instability comes from gradient vanishing。

# Fixing the gradient vanishing issue alone cannot stabilize training.

**Only Post-LN decoder suffers from gradient vanishing, but neither Post-LN Encoder, Pre-LN Encoder, nor Pre-LN Decoder.**

# Fixing the gradient vanishing issue alone cannot stabilize training.

**Only Post-LN decoder** Fix gradient vanishing **vanishing, but neither Post-LN Encoder, Pre-LN Encoder, nor Pre-LN Decoder.**

| Encoder | Decoder | Gradient | Training |
|---------|---------|----------|----------|
| Post-LN | Post-LN | Vanishing | Diverged |
| Post-LN | Pre-LN | ~~Vanishing~~ | Diverged |
| Pre-LN | Pre-LN | ~~Vanishing~~ | Converged |



Gradient vanishing only happens in backpropagations for Encoder-Attention sub-layers *i.e.,* from Encoder-Attention outputs to Self-Attention outputs.

11

# Unbalanced gradients are largely handled by adaptive optimizers.

| Relative Gradient Norm | $\dfrac{|\nabla\boldsymbol{w}_i^t|}{\max\limits_j|\nabla\boldsymbol{w}_j^t|}$ | Relative Parameter Update Norm | $\dfrac{|\boldsymbol{w}_i^{t+1}-\boldsymbol{w}_i^t|}{\max\limits_j|\boldsymbol{w}_j^{t+1}-\boldsymbol{w}_j^t|}$ |
|---|---|---|---|

As unbalanced gradients are largely handled by adaptive optimizers, it necessitates the use of adaptive optimizers.



Although the gradient distribution is unbalanced (*e.g.*, $W^{(V1)}$ and $W^{(V2)}$ have larger gradients than $W^{(K)}$ and $W^{(Q)}$), adaptive optimizers lead to consistent update magnitudes for different parameters.

Epoch # (iterations over the training set)

# Amplification Effect

Post-LN and Pre-LN aggregates residual branch outputs differently.

For a residual layer $x + f(x)$, we refer $f(x)$ as residual outputs and $x + f(x)$ as layer outputs

$\beta_{i,j}$ integrates all LNs and captures layer dependency.

Input of $(i+1)^{th}$ module/output of $i^{th}$ layer $= \sum_j \beta_{i,j} \cdot$ Normalized output of $j^{th}$ module

# $\beta_{i,j}$ integrates LNs and captures layer dependency



$\cdots\cdots\cdots\to = \beta_{3,3}\cdot\to\ +\beta_{3,2}\cdot\to\ +\beta_{3,2}\cdot$
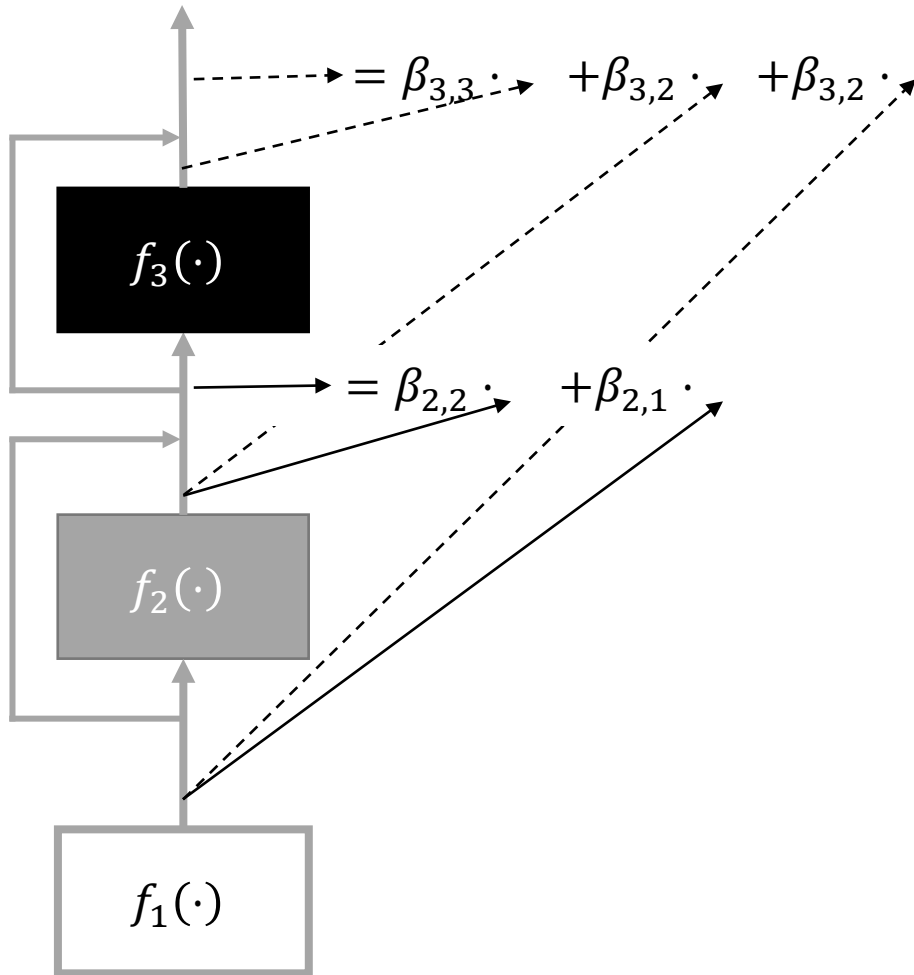
$f_3(\cdot)$

$= \beta_{2,2}\cdot\ +\beta_{2,1}\cdot$
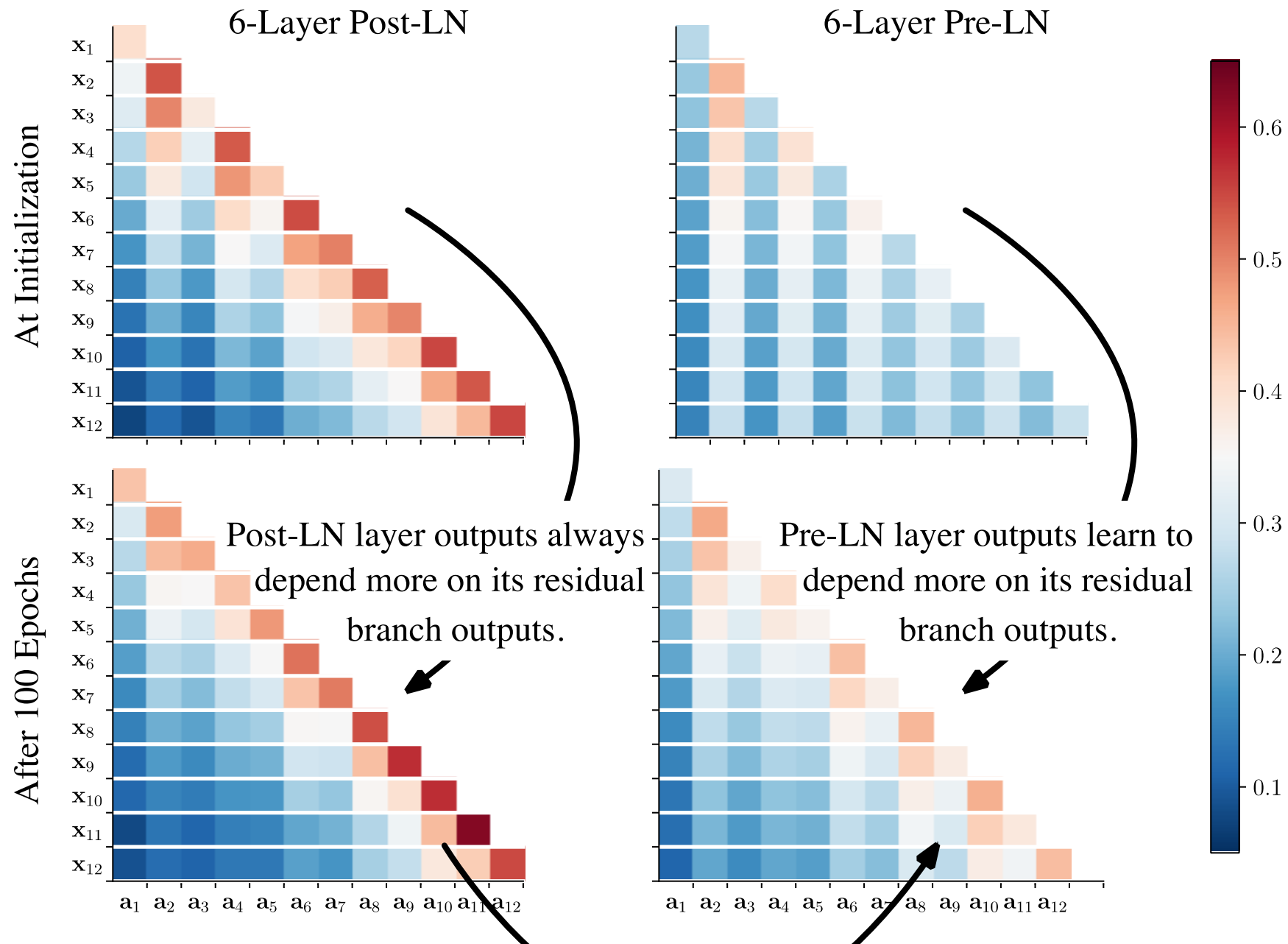
$f_2(\cdot)$

$f_1(\cdot)$

Refer $\beta_{i,i}$ as the dependency on its own residual branch.

standard deviation of $j^{\text{th}}$ output

For example, $\beta_{i,j} = \dfrac{\text{Std}[a_j]}{\text{Std}[\sum_{k \leq i} a_k]}$ for Pre-LN

standard deviation of the sum of the first i outputs.

Dependency on Residual Branches

6-Layer Post-LN

6-Layer Pre-LN

At Initialization

After 100 Epochs

Post-LN layer outputs always depend more on its residual branch outputs.

Pre-LN layer outputs learn to depend more on its residual branch outputs.

Comparing final models, Post-LN layer has a larger dependency on its residual branch.

# Large Dependency Destabilizes Training

Under some conditions, we have: $\text{Var}[\boxed{\mathcal{F}(\mathbf{x_0}, W) - \mathcal{F}(\mathbf{x_0}, W + \delta)}] \approx \sum_{i=1}^{N} \boxed{\beta_{i,i}^2} C$

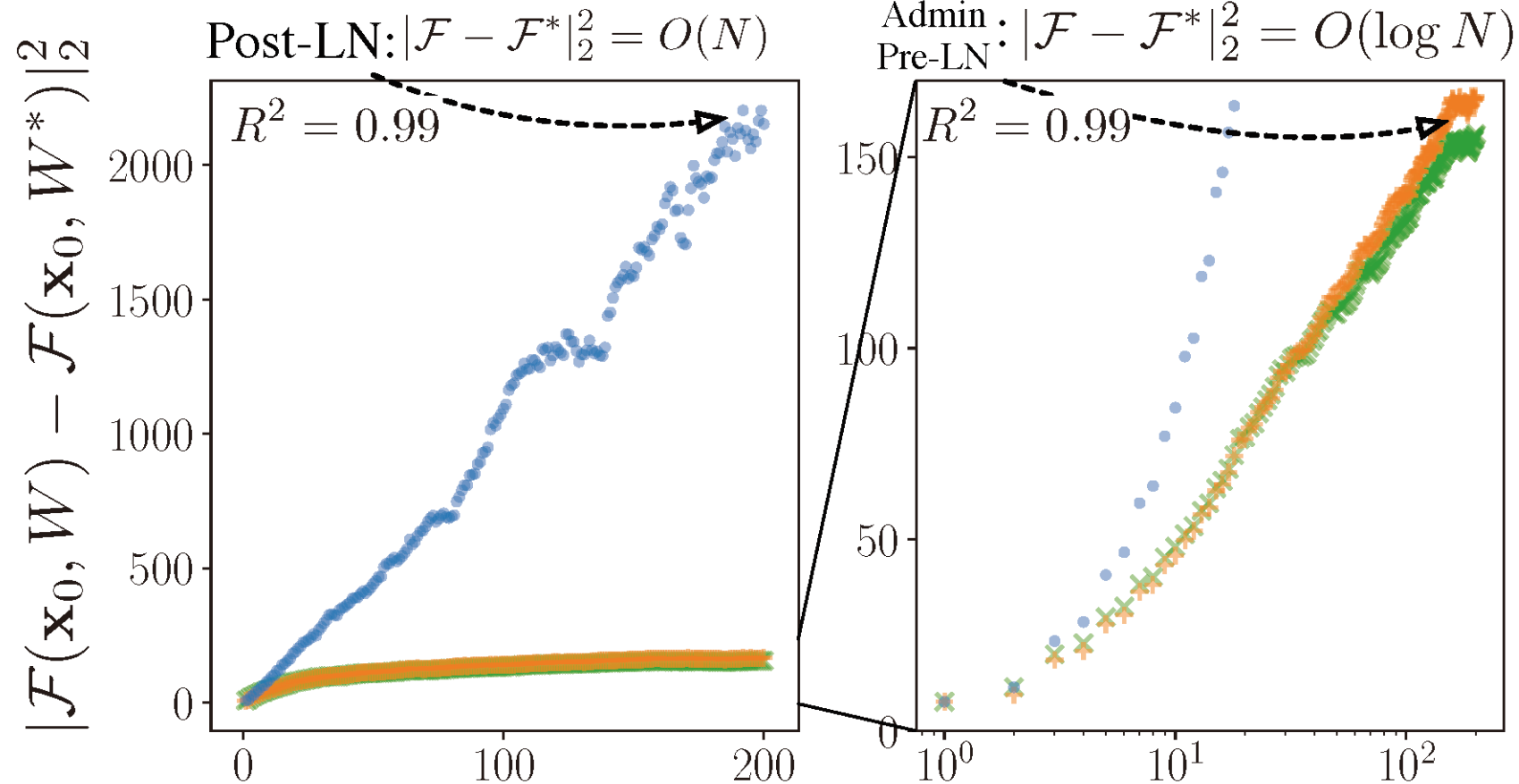Model output change.

Dependency on its own residual branch (the weight for i[th] residual outputs in i[th] layer outputs).

Corollary 1. For Pre-LN, $\text{Var}[\mathcal{F}(\mathbf{x_0}, W) - \mathcal{F}(\mathbf{x_0}, W + \delta)] = O(\log N)$ where N is layer #.

Corollary 2. For Post-LN, $\text{Var}[\mathcal{F}(\mathbf{x_0}, W) - \mathcal{F}(\mathbf{x_0}, W + \delta)] = O(N)$ where N is layer #.

# Large dependency destabilizes training



Random Perturbations, $i.e., W^* = W + \delta$

Post-LN: $|\mathcal{F} - \mathcal{F}^*|_2^2 = O(N)$

Admin
Pre-LN: $|\mathcal{F} - \mathcal{F}^*|_2^2 = O(\log N)$

Gradient Updates, $i.e.,$
$W^* = W + \mathrm{Adam}(\nabla_W \mathcal{L}(\mathcal{F}))$

$R^2 = 0.99$

$R^2 = 0.99$

$R^2 = 0.75$

Post-LN is less
stable than Pre-LN

Pre-LN
Post-LN
Admin (Post-LN)

$|\mathcal{F}(\mathbf{x}_0, W) - \mathcal{F}(\mathbf{x}_0, W^*)|_2^2$

Num of Sub-Layers (FFN or Self-Attention) in the Encoder

# Large dependency destabilizes training

Why *warmup* helps to alleviate the instability of Post-LN?

Under some conditions, we have: $\text{Var}[\mathcal{F}(\mathbf{x_0}, W) - \mathcal{F}(\mathbf{x_0}, W + \delta)] \approx \sum_{i=1}^{N} \beta_{i,i}^2 \boxed{C}$

Related to gradient norm

However, the difference between $O(\log N)$ and $O(N)$ is significant for deep networks (large N). In our experiments, <u>simply increasing the warmup steps</u> <u>fails</u> to stabilize the training of <u>deep Transformers</u> successfully.

# Model Initialization

We add $\boldsymbol{\omega_i}$ to restrict the layer dependency in the early stage of Post-LN.

$$\mathbf{x}_i = f_{\mathrm{LN}}(\mathbf{b_i}), \text{ where } \mathbf{b_i} = \mathbf{x_{i-1}} \cdot \boldsymbol{\omega_i} + f_i(\mathbf{x}_{i-1})$$

# Admin --- **Ad**aptive **m**odel **in**itialization

We add $\boldsymbol{\omega_i}$ to restrict the layer dependency in the early stage of Post-LN.

$$\mathbf{x}_i = f_{\mathrm{LN}}(\mathbf{b_i}), \text{where } \mathbf{b_i} = \mathbf{x_{i-1}} \cdot \boldsymbol{\omega_i} + f_i(\mathbf{x}_{i-1})$$

Also, we divide the initialization to two stages:

- Initialize $\boldsymbol{\omega_i}$ as 1, and empirically estimate the variance of $\mathrm{Var}[\mathbf{x}_i]$;

- Based on estimated variance, initialize $\boldsymbol{\omega_i}$ to ensure $\mathrm{Var}[\mathcal{F}(\mathbf{x_0}, W) - \mathcal{F}(\mathbf{x_0}, W + \delta)] = O(\log \mathrm{N})$ at initialization.

# Large dependency destabilizes training



Random Perturbations, $i.e., W^* = W + \delta$

Post-LN: $|\mathcal{F} - \mathcal{F}^*|_2^2 = O(N)$

Admin, Pre-LN: $|\mathcal{F} - \mathcal{F}^*|_2^2 = O(\log N)$

$R^2 = 0.99$

$R^2 = 0.99$

Gradient Updates, $i.e.,$

$W^* = W + \text{Adam}(\nabla_W \mathcal{L}(\mathcal{F}))$

$R^2 = 0.75$

Post-LN is less stable than Pre-LN

$|\mathcal{F}(\mathbf{x}_0, W) - \mathcal{F}(\mathbf{x}_0, W^*)|_2^2$

Num of Sub-Layers (FFN or Self-Attention) in the Encoder

+ Pre-LN
• Post-LN
× Admin (Post-LN)

# Admin --- **Ad**aptive **m**odel **in**itialization



18-Layer Admin (Post-LN)

18-Layer Pre-LN

Normalized Layer Outputs

At Initialization

Admin stabilizes model training by avoid over-large dependencies.

After 100 Epochs

With a Post-LN structure, Admin allows layer outputs of the final model to depend more on their residual branches.

Normalized Outputs of Residual Branches

# Experiments

# Experiments

| Dataset | IWSLT'14 De-En | WMT'14 En-Fr | | WMT'14 En-De | | |
|---|---|---|---|---|---|---|
| Enc #-Dec # | 6L-6L (small) | 6L-6L | 60L-12L | 6L-6L | 12L-12L | 18L-18L |
| Post-LN | 35.64±0.23 | 41.29 | failed | 27.80 | failed | failed |
| Pre-LN | 35.50±0.04 | 40.74 | 43.10 | 27.27 | 28.26 | 28.38 |
| Admin | **35.67±0.15** | **41.47** | **43.80** | **27.90** | **28.58** | **29.03** |

Without introducing any additional hyper-parameters, it achieves the new state-of-the-art on WMT'14 En-Fr (w.o. additional supervision including back translation).

# Experiments

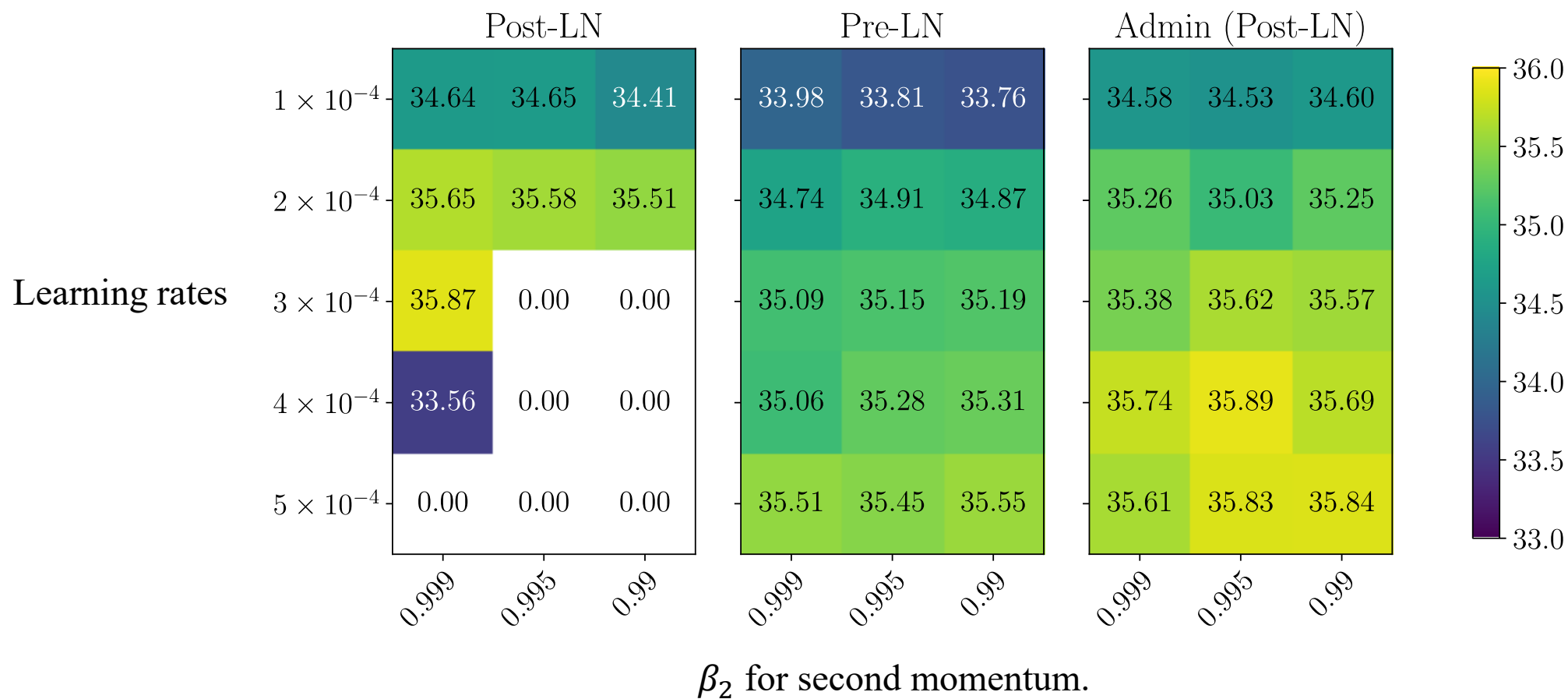| Dataset | IWSLT'14 De-En | WMT'14 En-Fr | | WMT'14 En-De | | |
|---------|----------------|--------------|--------------|--------------|--------------|--------------|
| Enc #-Dec # | 6L-6L (small) | 6L-6L | 60L-12L | 6L-6L | 12L-12L | 18L-18L |
| Post-LN | 35.64±0.23 | 41.29 | failed | 27.80 | failed | failed |
| Pre-LN | 35.50±0.04 | 40.74 | 43.10 | 27.27 | 28.26 | 28.38 |
| Admin | **35.67±0.15** | **41.47** | **43.80** | **27.90** | **28.58** | **29.03** |

We systematically evaluate deep Admin networks and summarizes results in the following report:

Liu, X., Duh, K., Liu, L., & Gao, J. (2020). Very deep transformers for neural machine translation. arXiv preprint arXiv:2008.07772.

Highlights: **30.1** BLEU on WMT'14 En-De, **46.4** BLEU on WMT'14 En-Fr (**w. back-translation**)

# Experiments



BLEU on IWSLT'14

# Take Away

- Unbalanced gradient is not the root cause of the unstable Transformer training. It is largely addressed by adaptive optimizers.

- Large dependencies on residual branches amplifies the fluctuation and destabilizes training.

- Controlling such dependencies at initialization, Admin is more stable, converges faster, and leads to better performance.