

Stable-Transformer: TOWARDS A STABLE TRANSFORMER TRAINING

Anonymous authors

Paper under double-blind review

ABSTRACT

The scale of parameters in Transformers has expanded dramatically—from hundreds of millions to several trillion. A key challenge when scaling the model to trillions is the training instability. Although many practical tricks, such as learning rate warmup, query-key normalization and better weight initialization, have been introduced to mitigate the training instability, a rigorous mathematical understanding of why such instabilities happen and why the above-mentioned tricks work well is still unclear. In this paper, we give a theoretical analysis of the initialization, normalization and attention mechanism in Transformers, and present a set of stabilized designs of the initialization, normalization and attention mechanism, which are thus termed as *StableInit*, *StableNorm* and *StableAtten*, individually. In experiments, we demonstrate that each of our stabilized designs, *i.e.*, *StableInit*, *StableNorm* and *StableAtten*, exhibits better stability. Furthermore, by putting the stabilized designs together, we propose a stabilized Transformer, termed *Stable-Transformer*, and show in experiments on large model (1B parameters) and deep model (200 layers) that our proposed *Stable-Transformer* achieves a more stable training process.

“My work always tried to unite the truth with the beautiful, but when I had to choose one or the other, I usually chose the beautiful.”

— Hermann Weyl

1 INTRODUCTION

The scale of parameters in Transformers (Vaswani et al., 2017; Radford et al., 2018; 2019; Brown et al., 2020; Touvron et al., 2023; Chowdhery et al., 2023) has expanded dramatically—from hundreds of millions to several trillion—parallel to significant advancements of hardware capabilities in the field of deep learning (Goodfellow et al., 2016; LeCun et al., 2015; Bengio et al., 2021). This exponential growth in model size has been facilitated by equally significant strides in computational power, enabling deeper and more complex network architectures. As these models have grown, they have set new benchmarks across a myriad of tasks in various fields such as natural language processing (Dubey et al., 2024; Achiam et al., 2023), computer vision (Ravi et al., 2024), and generation (Peebles & Xie, 2023).

Despite of these significant achievements, training larger models still suffers from an instability issue, which is often characterized as the difficulties in convergence, the sensitivity to initial conditions, and the necessity to finely tuned optimization strategies. Since that the instability in training process encumbers the deployment and real-world applicability of the sophisticated models, it is crucial to have a mathematical understanding of why such instability happens and it is urgent to invent stabilized design of the architecture or training strategies.

To gain a deeper understanding of the instability in training Transformers, it is essential to investigate the training dynamics of Transformers. To date, there are various studies devoted to the training process of Transformers from different perspectives, including normalization (Wang et al., 2019; Xiong et al., 2020; Liu et al., 2020; Miyato et al., 2018; Wang et al., 2022), attention mechanisms (Henry et al., 2020; Wortsman et al., 2024), model structures (Bachlechner et al., 2021; Qi et al., 2023b), and initialization strategies (Glorot & Bengio, 2010; He et al., 2015; Qi et al.,

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

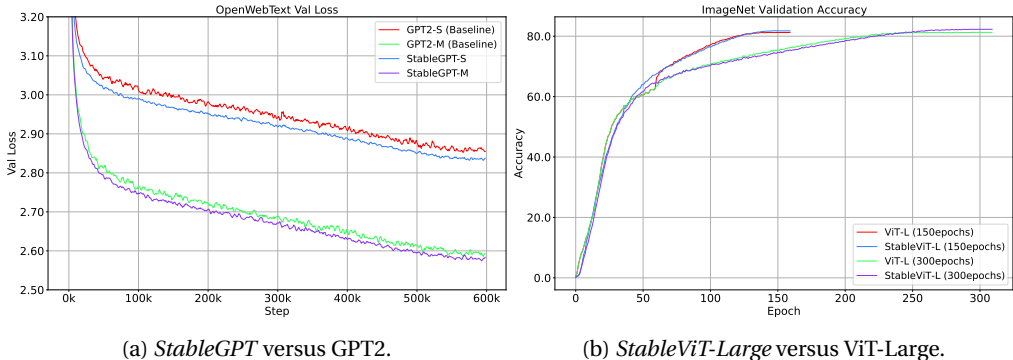


FIGURE 1: Except for being more stable during the training, *StableGPT* (left) also achieves a better validation loss (**2.827** for *StableGPT-S* (124M) versus **2.848** for GPT2-S (124M), and **2.569** for *StableGPT-M* (350M) versus **2.579** for GPT2-M (350M)), and *StableViT-L* (right) achieves a better recognition accuracy (**82.4% versus 81.3%**) compared to ViT-L under the settings of training 150 or 300 epochs. *StableGPT-S* can get an even better result (validation loss (**2.819**) with a higher learning rate (see Appendix I.), but here for fairness, we keep all the parameters of the optimizer in training are the same as the baseline.

2023b). From the perspective of normalization, it has shown that normalization play a critical role in stabilizing the training of Transformer, *e.g.*, Xiong et al. (2020) demonstrated that Pre-LayerNorm (Pre-LN) offers greater stability compared to Post-LayerNorm (Post-LN), Wang et al. (2022) proposed a DeepNorm and a depth-specific initialization to stabilize Post-LN. From the perspective of attention, Kim et al. (2021) showed that the standard dot-product attention is not Lipschitz continuous and thus introduced an alternative L_2 attention to address the continuous issue, QKNorm (Henry et al., 2020; Dehghani et al., 2023) proposed to normalize the query and key matrices in attention mechanisms to improve the stability of the attention module. From the perspective of initialization, Zhang et al. (2019) introduced a fixed-update initialization (Fixup) to prevent gradient exploding or vanishing at the start of training. This method rescales a standard initialization and enables stable training of residual networks without the need for normalization. Each of these approaches contributes to a better understanding and improvement of Transformer training stability, paving the way for more robust and efficient models. Bachlechner et al. (2021) demonstrated that a simple architectural modification, *i.e.*, gating each residual shortcut with a learnable zero-initialized parameter (ReZero), could significantly stabilize the training of Transformer. Using ReZero, they successfully trained Transformers with up to 120 layers. More recently, Qi et al. (2023b) introduced a novel Transformer architecture, called LipsFormer, which is designed to be Lipschitz continuous (*i.e.*, the gradients are bounded) and has been shown more stable during the training. The Lipschitz continuity allows for certain theoretical guarantees about the model’s behavior, which is important in reliability or interpretability.

This paper attempts to provide a theoretical understanding of the components that cause training instability of Transformer. To be specific, our main contributions are highlighted as follows.

- We give a theoretical analysis of the Xavier initialization from the perspective of random matrix theory, showing that the Lipschitz constant of the linear projection associated to the Xavier initialization is bounded by 2. Instead, we present a more stable method for initialization, termed *StableInit* (defined in Eq. 1), for which the Lipschitz constant is bounded by 1.
- We dig into the issue in back-propagation of the normalization by analyzing the Jacobian matrix of the normalization layer and find that the factor \sqrt{d} will affect the gradient flow significantly. As a remedy, we derive a more stable design for the normalization, called *StableNorm* (defined in Eq. 2), in which d^α with $\alpha \in [0, 0.5]$ is adopted to replace \sqrt{d} in the normalization layer, and verify that using a smaller α (*e.g.*, 0.475, rather than 0.5) will yield smaller gradients and thus lead to more stable training.
- We present a new stable form of attention, named *StableAtten* (defined in Eq. 3), which is built on our *StableNorm* and has the advantage that the logit of the attention is not directly related to the hidden dimension d and thus is robust to the increase of the model scale.

- By putting together the *StableInit*, *StableNorm* and *StableAtten*, we have a stabilized design for Transformer, termed *Stable-Transformer*. In experiments, for the single-direction generative model *i.e.*, GPT (Radford et al., 2018; 2019; Brown et al., 2020), we compile a *StableGPT*; for a bi-direction attention model *i.e.*, ViT (Dosovitskiy et al., 2020), we compile a *StableViT*. We evaluate *StableGPT* and *StableViT* extensively on large model (1B parameters) and deep model (200 layers) that our proposed *Stable-Transformer* achieves a more stable training process.

The paper is organized as follows. We first introduce our experimental setups, and then present our stabilized design of the modules. For each module, we give our mathematical analysis at first and then show empirical evaluation. There is no an independent section for experiments.

2 EXPERIMENTAL METHODOLOGY

We evaluate our stabilized components on ViT (Dosovitskiy et al., 2020) and GPT (Radford et al., 2018; 2019; Brown et al., 2020). For general training setting, by default, we use the optimizer Adam Kingma & Ba (2014) with $\beta_1 = 0.9$, $\beta_2 = 0.95$ and $\epsilon = 10^{-8}$, and the gradient clipping is set to 1. When using weight decay, we follow AdamW (Loshchilov & Hutter, 2019), for which only the weight matrix is enforced to the weight matrix but not the 1-d vector (*e.g.*, γ and β in LayerNorm) and the scalar. We train all models on GPUs A800 in bfloat16 precision using PyTorch (Paszke et al., 2019). We use a cosine-decay (Loshchilov & Hutter, 2016) schedule from a preset maximum learning rate to a preset minimum learning rate.

Experimental setups for ViT (Dosovitskiy et al., 2020) and our *StableViT*. We use timm Wightman (2019)¹. For ViT model, we use two different scales: ViT-Large (ViT-G) and ViT-Huge (ViT-H). The detailed information about these models are summarized in Table 1. For data augmentation, we use the same data augmentation as Adan (Xie et al., 2024). Thus our results are aligned with the results reported in (Xie et al., 2024).

Experimental setups for GPT2 and our *StableGPT*. We use nanoGPT² (Karpathy, 2022), which is a simple and fast repository for training and fine-tuning the medium-sized GPTs. The GPT2 is implemented in four versions: GPT2-Small (GPT2-S), GPT2-Medium (GPT2-M), GPT2-Large (GPT2-L) and GPT2-XL. Due to time and computational costs, we only use GPT2-Small (GPT2-S), GPT2-Medium (GPT2-M).

We align our experiments with the original repository, and use the exactly same training setting as nanoGPT. Detailed parameters is listed in Table 3. We reproduce the baseline GPT-2 124M model with the same setup as in nanoGPT, for which the training loss is reduced to 2.848. The learning curve of the loss matches to the original nanoGPT.

It is worth to note that when evaluating a module (or method), we keep all the same but the specific module (or method) for fair comparison. To be more specific, when evaluating each of *StableInit*, *StableNorm* and *StableAtten*, we only replace the corresponding module (or method).

3 STABLE-TRANSFORMER AND ITS THEORETICAL JUSTIFICATIONS

In this section, we will present our stabilized initialization method *StableInit*, stabilized normalization module *StableNorm*, and stabilized attention mechanism *StableAtten*, individually. Moreover, we will combine them together to build our *Stable-Transformer*. For each method or module, we start with a justification for the instability issue in training and then provide our stabilized designs with both theoretical justification and empirical evaluation.

3.1 STABLE INITIALIZATION

The Xavier initialization is a remarkable technique that significantly enhances the training of neural networks by initializing the weights in a way that maintains the variance of activation across layers, to mitigate the vanishing or exploding gradient problem. In Glorot & Bengio (2010), the Xavier initialization is sorted to two types, *i.e.*, Gaussian distribution and uniform

¹<https://github.com/huggingface/pytorch-image-models/tree/main>

²<https://github.com/karpathy/nanoGPT>

distribution. The Xavier initialization for $\mathbf{W} \in \mathbb{R}^{n_{in} \times n_{out}}$ with Gaussian distribution is defined as: $W_{i,j} \stackrel{i.i.d.}{\sim} \mathcal{N}\left(0, \frac{2}{n_{in} + n_{out}}\right)$, where n_{in} and n_{out} denote the dimensions of the input and the output. It is widely used in training modern neural networks and is usually as the default initialization method. Therefore, in the following we will only consider the Xavier initialization with Gaussian distribution when mentioning it.

Now we will analyze the property of the Xavier initialization with Gaussian distribution. To begin with, we would like to introduce a theorem from Random Matrix Theory (RMT) (Wigner, 1955; Tao, 2012; Edelman & Rao, 2005). From RMT, we have the theorem about the singular values of a Gaussian random matrix.

Theorem 1 (Singular Value Bounds of a Gaussian Random Matrix)

Let $\mathbf{W} \in \mathbb{R}^{m \times n}$ have i.i.d. standard Gaussian entries, i.e., $W_{i,j} \stackrel{iid}{\sim} \mathcal{N}(0, 1)$. For every $m \geq n$, we have the following inequality

$$\sqrt{m} - \sqrt{n} \leq \mathbb{E}[\sigma_{\min}(\mathbf{W})] \leq \mathbb{E}[\sigma_{\max}(\mathbf{W})] \leq \sqrt{m} + \sqrt{n},$$

where $\sigma_{\min}(\mathbf{W})$ and $\sigma_{\max}(\mathbf{W})$ denote the minimal and maximal singular values, respectively.

We provide the proof of Theorem 1 via theory from high-dimensional probability (Vershynin, 2010) in the appendix B. Theorem 1 presents that a random matrix initialized by standard Gaussian distribution $\mathcal{N}(0, 1)$, the expectations of its largest and smallest singular values are bounded. The expectation of its largest singular value is no more than $\sqrt{m} + \sqrt{n}$, and the expectation of its smallest singular value is no less than $\sqrt{m} - \sqrt{n}$.

According to Theorem 1, and the definition of Xavier initialization, we have the following lemma.

Lemma 1 (Upper Bound of Weight Matrix Norm of Xavier Initialization)

Let $\mathbf{W} \in \mathbb{R}^{n_{in} \times n_{out}}$ have i.i.d. standard Gaussian entries, i.e., $W_{i,j} \stackrel{iid}{\sim} \mathcal{N}\left(0, \frac{2}{n_{in} + n_{out}}\right)$. we have the following inequality for its maximum singular value, $\mathbb{E}[\sigma_{\max}(\mathbf{W})] \leq 2$.

We provide a proof of Lemma 1 in Appendix E.

Remark 1. Back to the year 2010 when we still do not have ResNet (He et al., 2016) and Batch Normalization (Ioffe & Szegedy, 2015), the Xavier initialization is a remarkable technique that enables the researchers to train a network more than 10 layers. Suppose that we have a MLP with 10 linear layers with ReLU (Nair & Hinton, 2010) between two nearby linear layers, with a softmax in the last linear layer, and using a cross-entropy loss, then the expectation of the largest singular value of each layer is up to 2 and it means that the Lipschitz constant (Fazlyab et al., 2019; Kim et al., 2021; Qi et al., 2023a;b) for each linear layer is 2. Therefore, we can compute the Lipschitz constant of the whole network (assuming that the softmax and the cross-entropy has Lipschitz constant 1) as $2^{10} = 1024$, which is under control.

Although Xavier initialization is a popular initialization method, it still has some issues. One main disadvantage in Xavier initialization is that *it is sensitive to the increase of the network depth*. If the number of layers is 50, then the Lipschitz constant of the above-mentioned whole network will be extremely huge. To fix this issue, we present a simple and 1-Lipschitz initialization strategy, which is termed **StableInit**. Precisely, we define it as follows:

$$W_{i,j} \stackrel{iid}{\sim} \mathcal{N}\left(0, \left(\frac{1}{\sqrt{n_{in}} + \sqrt{n_{out}}}\right)^2\right), \quad (1)$$

where n_{in} and n_{out} are the dimensions of the input and output of the module, respectively. It is easy to see that with our *StableInit* initialization, we have that $\mathbb{E}[\sigma_{\max}(\mathbf{W})] \leq 1$. Similar to the Xavier initialization, we can also consider to multiply a gain on the weight. The gain can be set to be a smaller value when we use deeper networks. By default, we use 1.0 as the gain.

For clarity, we summarize the following two properties of our *StableInit* initialization.

- **1-Lipschitz constant:** with our *StableInit*, a linear projection module has a Lipschitz constant with expectation 1, other than 2 in the original Xavier initialization.

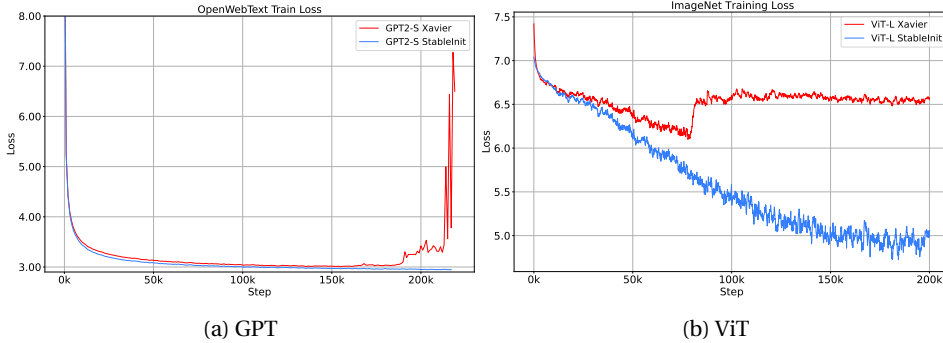


FIGURE 2: Evaluation of *StableInit* and comparing to the original Xavier initialization. In the legend, “GPT2-S” denotes GPT2-Small. To compare the stability, we do not use learning rate warmup in the evaluation.

- **Less sensitive to depth increase:** our *StableInit* is less sensitive to the depth increase compared to the the original initialization. For a MLP with 10 or 100 linear layers with a ReLU (Nair & Hinton, 2010) between two linear layers, the Lipschitz constant is 1 under our *StableInit*.

3.1.1 EVALUATION FOR *StableInit*

These properties of our stabilized Xavier initialization will lead to a more stable training compared to the original Xavier initialization because it has a proper Lipschitz constant.

We can see that from Figure 2 *StableInit* obtains a more stable training property. With *StableInit*, the model can be trained longer until diverge. Theoretically, *StableInit* will be more robust than Xavier initialization for larger model. In experiments, the learning curve of ViT is more jitter because its supervision signal is more sparse, a batch of tokens of GPT is 0.5M tokens, each token will provide a signal, but the batch size 1024 in ViT only provides 1024 supervision signals.

3.2 STABLE NORMALIZATION

LayerNorm (Ba et al., 2016) is a technique widely used in deep learning to stabilize and accelerate the training of neural networks. The original definition of LayerNorm is $\text{LN}(x) = \gamma \odot z + \beta$, where $z = \frac{y}{\text{std}(y)}$ and $y = (\mathbf{I} - \frac{1}{d}\mathbf{1}\mathbf{1}^\top)x$. After adding a smoothing factor, it can also be written as, $\text{LN}(x) = \gamma \odot \frac{\sqrt{d}y}{\sqrt{\|y\|_2^2 + \epsilon}} + \beta$, and $y = (\mathbf{I} - \frac{1}{d}\mathbf{1}\mathbf{1}^\top)x$, where ϵ is the smoothing factor, d is the feature dimension of x , γ and β are two learnable \mathbb{R}^d vectors, γ and β are initialized to 1 and 0. Most recently, some new large language models (Touvron et al., 2023; Chowdhery et al., 2023; Team, 2023) uses RMSNorm (Zhang & Sennrich, 2019) to replace LayerNorm, where RMSNorm is defined as: $\text{RMSN}(x) = \gamma \odot \frac{\sqrt{d}x}{\sqrt{\|x\|_2^2 + \epsilon}}$. Compared to LayerNorm, RMSNorm does not use the bias term and does not conduct the centering.

The Jacobian matrices of LayerNorm and RMSNorm with respect to x are calculated as follows:

$$\frac{\partial \text{LN}(x)}{\partial x} = \frac{\partial y}{\partial x} \frac{\partial \text{LN}(x)}{\partial y} = \frac{\sqrt{d}}{\sqrt{\|y\|_2^2 + \epsilon}} \left(\mathbf{I} - \frac{1}{d}\mathbf{1}\mathbf{1}^\top \right) \left(\mathbf{I} - \frac{yy^\top}{\|y\|_2^2 + \epsilon} \right) \text{diag}(\gamma),$$

$$\frac{\partial \text{RMSN}(x)}{\partial x} = \frac{\sqrt{d}}{\sqrt{\|x\|_2^2 + \epsilon}} \left(\mathbf{I} - \frac{xx^\top}{\|x\|_2^2 + \epsilon} \right) \text{diag}(\gamma).$$

Let us explain each term a little bit individually. It is easy to prove that the maximal singular value of $(\mathbf{I} - \frac{1}{d}\mathbf{1}\mathbf{1}^\top)$ and $(\mathbf{I} - \frac{yy^\top}{\|y\|_2^2 + \epsilon})$ are both 1. We give a proof of $\sigma_{\max}(\mathbf{I} - \frac{yy^\top}{\|y\|_2^2 + \epsilon}) \leq 1$ in Appendix C. Note that $\sigma_{\max}(\mathbf{I} - \frac{1}{d}\mathbf{1}\mathbf{1}^\top) \leq 1$ is a special case of the former.

To analyze and compare $\frac{\sqrt{d}}{\sqrt{\|\mathbf{y}\|_2^2 + \epsilon}}$ in $\frac{\partial \text{LN}(\mathbf{x})}{\partial \mathbf{x}}$ and $\frac{\sqrt{d}}{\sqrt{\|\mathbf{x}\|_2^2 + \epsilon}}$ in $\frac{\partial \text{RMSN}(\mathbf{x})}{\partial \mathbf{x}}$, we have the following inequality for the centering transformation.

Theorem 2 (Centering Transformation Inequality)

Let $\mathbf{y} = (\mathbf{I} - \frac{1}{d} \mathbf{1} \mathbf{1}^\top) \mathbf{x}$, we have the following inequality: $\frac{\sqrt{d}}{\sqrt{\|\mathbf{x}\|_2^2 + \epsilon}} \leq \frac{\sqrt{d}}{\sqrt{\|\mathbf{y}\|_2^2 + \epsilon}}$.

Proof. Centering can be denoted as, $\mathbf{y} = (\mathbf{I} - \frac{1}{d} \mathbf{1} \mathbf{1}^\top) \mathbf{x} = \mathbf{x} - \frac{1}{d} (\sum_{i=1}^d x_i) \mathbf{1}$. Then we have,

$$\|\mathbf{y}\|_2^2 = \left(\mathbf{x} - \frac{1}{d} \left(\sum_{i=1}^d x_i \right) \mathbf{1} \right)^\top \left(\mathbf{x} - \frac{1}{d} \left(\sum_{i=1}^d x_i \right) \mathbf{1} \right) = \mathbf{x}^\top \mathbf{x} - 2 \frac{1}{d} \left(\sum_{i=1}^d x_i \right) \mathbf{1}^\top \mathbf{x} + \frac{1}{d^2} \left(\sum_{i=1}^d x_i \right)^2 \mathbf{1}^\top \mathbf{1}$$

Since $\mathbf{1}^\top \mathbf{x} = \sum_{i=1}^d x_i$ and $\mathbf{1}^\top \mathbf{1} = d$, we get: $\|\mathbf{y}\|_2^2 = \|\mathbf{x}\|_2^2 - \frac{1}{d} (\sum_{i=1}^d x_i)^2$. Since the term $\frac{1}{d} (\sum_{i=1}^d x_i)^2$ is non-negative, we have $\|\mathbf{y}\|_2^2 \leq \|\mathbf{x}\|_2^2$. Therefore, we have $\frac{\sqrt{d}}{\sqrt{\|\mathbf{x}\|_2^2 + \epsilon}} \leq \frac{\sqrt{d}}{\sqrt{\|\mathbf{y}\|_2^2 + \epsilon}}$, which proves the inequality. \square

We can see that $\frac{\sqrt{d}}{\sqrt{\|\mathbf{y}\|_2^2 + \epsilon}}$ in $\frac{\partial \text{LN}(\mathbf{x})}{\partial \mathbf{x}}$ reaches to the maximum value $\frac{\sqrt{d}}{\sqrt{\epsilon}}$ when $\text{std}(\mathbf{x})$ is 0, but $\frac{\sqrt{d}}{\sqrt{\|\mathbf{x}\|_2^2 + \epsilon}}$ in $\frac{\partial \text{RMSN}(\mathbf{x})}{\partial \mathbf{x}}$ reaches to 0 if and only if \mathbf{x} is equal to $\mathbf{0}$. Theorem 2 means that RMSNorm is less likely to obtain the maximum value compared to LayerNorm.

We also observe that in the Jacobian matrices of LayerNorm and RMSNorm, both of them have a term \sqrt{d} , which is the dimension of \mathbf{x} and is also called the hidden dimension of the networks in large language model. With the increase of the hidden dimension d in larger models, there is a square root ratio effect, and thus may lead to larger gradients. Therefore, it will make it harder to train larger models. To alleviate this issue, we present a simple but more stable normalization mechanism as follows:

$$\text{StableNorm}(\mathbf{x}) = \gamma \odot \frac{d^\alpha \mathbf{x}}{\sqrt{\|\mathbf{x}\|_2^2 + \epsilon}}, \quad (2)$$

where α is a hyper-parameter, the range of α is $[0, 0.5]$. By choosing a reasonable α , we can obtain a more stable normalization module. We term our stabilized normalization as **StableNorm**. When $\alpha = 0.5$, **StableNorm** will be equal to RMSNorm (Zhang & Sennrich, 2019).

It is easy to derive that the Jacobian matrix of **StableNorm** is

$$\frac{\partial \text{StableNorm}(\mathbf{x})}{\partial \mathbf{x}} = \frac{d^\alpha}{\sqrt{\|\mathbf{x}\|_2^2 + \epsilon}} \left(\mathbf{I} - \frac{\mathbf{x} \mathbf{x}^\top}{\|\mathbf{x}\|_2^2 + \epsilon} \right) \text{diag}(\gamma).$$

We can see that the maximum value is $\frac{d^\alpha}{\sqrt{\epsilon}}$. A reasonable choice for ϵ is 10^{-5} . Along with the increase of the hidden dimension in larger model, we can tune the α to make the normalization layer more stable. A good strategy is to choose a smaller α for larger model. To have a more intuitive understanding, let us see an example. Assume $d = 4096$, then $d^\alpha = 64$ when $\alpha = 0.5$ whereas $d^\alpha \approx 42.22$ when $\alpha = 0.45$, it means that the gradient will be scaled by $4096^{0.45} \approx 42.22$ in our **StableNorm** instead of $4096^{0.5} = 64$. Similarly, when $\alpha = 0.4$, the gradient will be scaled by $4096^{0.4} \approx 27.85$ instead of 64 with $\alpha = 0.5$. By choosing 0.45 instead of 0.5, we can actually scale down the gradient by a factor of $\frac{42.22}{64} = 0.66$. Therefore, by choosing a reasonable α , our **StableNorm** will help to yield more stable gradients compared to RMSNorm and LayerNorm.

3.2.1 EVALUATION FOR **StableNorm**

According to our above derivation, different choices of α in Eq. (2) lead to different Jacobian matrix, and thus lead to different gradient flow. Here, we conduct a set of evaluations with different choices of α . We evaluate five different choices of α , i.e., $\alpha \in \{0.0, 0.125, 0.25, 0.375, 0.5\}$. We conduct experiments on GPT and ViT, respectively. The empirical results are shown in Figure 3.

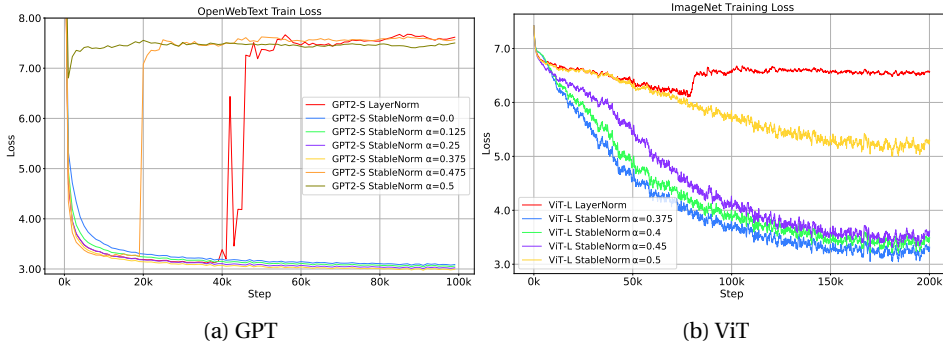


FIGURE 3: Evaluation of *StableNorm* and comparison to the original LayerNorm and RMSNorm on ViT and GPT, respectively. For instance, “StalbeGPT-S $\alpha = 0.5$ ” in the legend denotes StableGPT-Small model with α set to be 0.5. Where using $\alpha = 0.5$ is reduced to RMSNorm. To compare the stability, we do not use learning rate warmup.

we can see that from Figure 3, choosing an extremely small α may lead to gradient vanishing issue, but choosing a large α may causing training instability. How to choose a good α is an empirical trick. We note here that some relatively large α can be selected for GPT and some relatively small α can be selected for ViT, which may be related to the density and sparsity of the supervised signal of GPT and ViT. When d is large in some large model, a reasonable choice is to choose a smaller α (which takes a lot of resources to verify this. We thus do not do it in this paper). According to our current experiments, we can find that α is a good parameter for balancing gradient vanishing and exploding.

3.3 STABLE ATTENTION

Before we present our stabilized attention module, we review the self-attention (Vaswani et al., 2017) and the self-attention with Query-Key normalization (QKNorm) (Henry et al., 2020; Dehghani et al., 2023) at first. Rather than showing the QKNorm works as in Wortsman et al. (2024), we focus on revealing the underlying theoretical reasons. Then we will present our stabilized module, termed *StableAtten*.

3.3.1 WHY SELF-ATTENTION WITH QKNORM WORKS?

For an input sequence $\mathbf{X} \in \mathcal{R}^{d \times l}$, d is the dimension of the feature and l is the sequence length, self-attention (Vaswani et al., 2017) is defined as:

$$\mathbf{Y} = \mathbf{W}_v \mathbf{X} \mathbf{A}, \quad \mathbf{A} = \text{softmax}(\mathbf{P}^{(1)}), \quad \text{where } \mathbf{P}^{(1)} = \frac{\mathbf{X}^\top \mathbf{W}_q^\top \mathbf{W}_k \mathbf{X}}{\sqrt{d_1}},$$

in which $d_1 = d/h$ and h is the number of heads, d_1 is called as the head dimension, $\mathbf{W}_q, \mathbf{W}_k \in \mathcal{R}^{d_1 \times d}$, $\mathbf{W}_v \in \mathcal{R}^{d_2 \times d}$ (in practice, we usually set $d_1 = d_2$). The size of the output $\mathbf{Y} \in \mathcal{R}^{d_2 \times l}$, and the attention matrix $\mathbf{A} \in \mathcal{R}^{l \times l}$. Therefore, we have $P_{ij}^{(1)} = \frac{\mathbf{x}_i^\top \mathbf{W}_q^\top \mathbf{W}_k \mathbf{x}_j}{\sqrt{d_1}}$ where $\mathbf{P}^{(1)}$ is called the logit in (Dehghani et al., 2023; Wortsman et al., 2024).

Self-attention with QKNorm (Dehghani et al., 2023) is defined as:

$$\mathbf{A} = \text{softmax}(\mathbf{P}^{(2)}), \quad \text{where } P_{ij}^{(2)} = \left(\frac{\text{RMSN}(\mathbf{W}_q \mathbf{x}_i)^\top \text{RMSN}(\mathbf{W}_k \mathbf{x}_j)}{\sqrt{d_1}} \right).$$

Note that Dehghani et al. (2023) use a LayerNorm layer without bias and centering, which is equivalent to a RMSNorm Zhang & Sennrich (2019) layer as defined above.

For clarity, we write down the formula for the i -th query \mathbf{q}_i and the j -th key \mathbf{k}_j in QKNorm (Dehghani et al., 2023) as follows:

$$\begin{aligned}\mathbf{q}_i &= \text{RMSN}(\mathbf{W}_q \mathbf{x}_i) = \gamma_q \odot \frac{\sqrt{d_1} \mathbf{W}_q \mathbf{x}_i}{\sqrt{\|\mathbf{W}_q \mathbf{x}_i\|_2^2 + \epsilon}} = \sqrt{d_1} \text{diag}(\gamma_q) \frac{\mathbf{W}_q \mathbf{x}_i}{\sqrt{\|\mathbf{W}_q \mathbf{x}_i\|_2^2 + \epsilon}}, \\ \mathbf{k}_j &= \text{RMSN}(\mathbf{W}_k \mathbf{x}_j) = \gamma_k \odot \frac{\sqrt{d_1} \mathbf{W}_k \mathbf{x}_j}{\sqrt{\|\mathbf{W}_k \mathbf{x}_j\|_2^2 + \epsilon}} = \sqrt{d_1} \text{diag}(\gamma_k) \frac{\mathbf{W}_k \mathbf{x}_j}{\sqrt{\|\mathbf{W}_k \mathbf{x}_j\|_2^2 + \epsilon}}.\end{aligned}$$

It is easy to see that

$$\begin{aligned}P_{ij}^{(2)} &= \frac{\mathbf{q}_i^\top \mathbf{k}_j}{\sqrt{d_1}} = \frac{\sqrt{d_1} \sqrt{d_1} \left(\frac{\mathbf{W}_q \mathbf{x}_i}{\sqrt{\|\mathbf{W}_q \mathbf{x}_i\|_2^2 + \epsilon}} \right)^\top (\text{diag}(\gamma_q))^\top \text{diag}(\gamma_k) \frac{\mathbf{W}_k \mathbf{x}_j}{\sqrt{\|\mathbf{W}_k \mathbf{x}_j\|_2^2 + \epsilon}}}{\sqrt{d_1}} \\ &= \sqrt{d_1} \left(\frac{\mathbf{W}_q \mathbf{x}_i}{\sqrt{\|\mathbf{W}_q \mathbf{x}_i\|_2^2 + \epsilon}} \right)^\top \text{diag}(\gamma_q) \text{diag}(\gamma_k) \frac{\mathbf{W}_k \mathbf{x}_j}{\sqrt{\|\mathbf{W}_k \mathbf{x}_j\|_2^2 + \epsilon}}.\end{aligned}$$

We find that when using QKNorm in self-attention, both the logit and the gradients are *upper bounded*. Precisely, we have the following theorem.

Theorem 3 (Logit Inequality for Self-attention with QKNorm)

The logit in self-attention with QKNorm, where γ_q and γ_k are initialized to 1 and not learned, is upper bounded, i.e., $|P_{ij}^{(2)}| < \sqrt{d_1}$.

Proof. We have that when γ_q and γ_k are initialized to 1, then $P_{ij}^{(2)} = \sqrt{d_1} \left(\frac{\mathbf{W}_q \mathbf{x}_i}{\sqrt{\|\mathbf{W}_q \mathbf{x}_i\|_2^2 + \epsilon}} \right)^\top \frac{\mathbf{W}_k \mathbf{x}_j}{\sqrt{\|\mathbf{W}_k \mathbf{x}_j\|_2^2 + \epsilon}} = \sqrt{d_1} \cos(\theta) \leq \sqrt{d_1}$, where θ is the angle between $\frac{\mathbf{W}_q \mathbf{x}_i}{\sqrt{\|\mathbf{W}_q \mathbf{x}_i\|_2^2 + \epsilon}}$ and $\frac{\mathbf{W}_k \mathbf{x}_j}{\sqrt{\|\mathbf{W}_k \mathbf{x}_j\|_2^2 + \epsilon}}$. \square

On contrary, the logit $P_{ij}^{(1)}$ in the original self-attention is not bounded, since that $P_{ij}^{(1)} = \frac{\mathbf{x}_i^\top \mathbf{W}_q^\top \mathbf{W}_k \mathbf{x}_j}{d_1}$, where \mathbf{x}_i , \mathbf{x}_j and \mathbf{W}_q or \mathbf{W}_k might be not bounded. The *upper bounded logit in self-attention with QKNorm* is one of the theoretical reasons that QKNorm leads stable training process. Under a fixed hidden dimension d , we see that increasing the number of heads h will correspond to decrease the head dimension d_1 . According to the upper bound, we have that a smaller d_1 will stabilize the training process. To verify it, we evaluate the influence of the number of heads in experiment, and show the results in Figure 4. We can observe that increasing the number of heads in the attention does stabilize the training process of the original Transformer.

On the other hand, we find that the gradients in self-attention with QKNorm is also upper bounded. Note that all modules are updated by error back-propagation (Rumelhart et al., 1986; LeCun et al., 2002; 1989; 1998), and computing the gradients in chain is the key. The gradients in self-attention without or with QKNorm can be calculated as follows:

$$\frac{\partial P_{ij}^{(1)}}{\partial \mathbf{x}_i} = \mathbf{W}_q^\top \mathbf{W}_k \mathbf{x}_j, \quad \frac{\partial P_{ij}^{(1)}}{\partial \mathbf{x}_j} = \mathbf{W}_k^\top \mathbf{W}_q \mathbf{x}_i, \quad \frac{\partial P_{ij}^{(1)}}{\partial \mathbf{W}_q} = \mathbf{x}_i \mathbf{x}_j^\top \mathbf{W}_k^\top, \quad \frac{\partial P_{ij}^{(1)}}{\partial \mathbf{W}_k} = \mathbf{W}_q^\top \mathbf{x}_i \mathbf{x}_j^\top.$$

and

$$\begin{aligned}\frac{\partial P_{ij}^{(2)}}{\partial \mathbf{x}_i} &= \frac{\sqrt{d_1}}{\sqrt{\|\mathbf{W}_q \mathbf{x}_i\|_2^2 + \epsilon}} \mathbf{W}_q^\top \left(\mathbf{I} - \frac{\mathbf{W}_q \mathbf{x}_i (\mathbf{W}_q \mathbf{x}_i)^\top}{\|\mathbf{W}_q \mathbf{x}_i\|_2^2 + \epsilon} \right) \text{diag}(\gamma_q) \text{diag}(\gamma_k) \frac{\mathbf{W}_k \mathbf{x}_j}{\sqrt{\|\mathbf{W}_k \mathbf{x}_j\|_2^2 + \epsilon}}, \\ \frac{\partial P_{ij}^{(2)}}{\partial \mathbf{W}_q} &= \frac{\sqrt{d_1}}{\sqrt{\|\mathbf{W}_q \mathbf{x}_i\|_2^2 + \epsilon}} \left(\mathbf{I} - \frac{\mathbf{W}_q \mathbf{x}_i (\mathbf{W}_q \mathbf{x}_i)^\top}{\|\mathbf{W}_q \mathbf{x}_i\|_2^2 + \epsilon} \right) \text{diag}(\gamma_q) \text{diag}(\gamma_k) \frac{\mathbf{W}_k \mathbf{x}_j}{\sqrt{\|\mathbf{W}_k \mathbf{x}_j\|_2^2 + \epsilon}} \mathbf{x}_i^\top.\end{aligned}$$

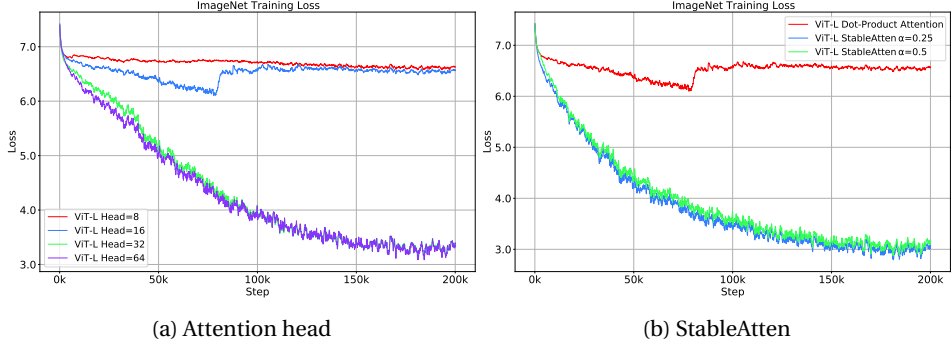


FIGURE 4: Evaluation of the number of attention head, and comparison of *StableAtten* with the original dot-product self-attention in Transformer using ViT. For instance, “ViT-L Head=8” in the legend denotes ViT-Large model with 8 attention heads. To compare the training stability, we do not use learning rate warmup in the evaluation.

Let us have a comparison between $\frac{\partial P_{ij}^{(1)}}{\partial \mathbf{x}_i}$ and $\frac{\partial P_{ij}^{(2)}}{\partial \mathbf{x}_i}$. We have the following observations.

- $\frac{\partial P_{ij}^{(1)}}{\partial \mathbf{x}_i}$ is not bounded, but $\frac{\partial P_{ij}^{(2)}}{\partial \mathbf{x}_i}$ is bounded by $\frac{\sqrt{d} \|\mathbf{W}_q\|}{\sqrt{\epsilon}}$. The upper bounded gradients lead to more stable training;
- The value of $\frac{\partial P_{ij}^{(1)}}{\partial \mathbf{x}_i}$ is proportion to $\mathcal{O}(\|\mathbf{W}_q^\top \mathbf{W}_k\|)$, but $\frac{\partial P_{ij}^{(2)}}{\partial \mathbf{x}_i}$ is only proportion to $\mathcal{O}(\|\mathbf{W}_q\|)$. Generally speaking, the spectral norm of $\mathcal{O}(\|\mathbf{W}\|)$ will increase along with the training process and will saturate and oscillate when training comes to convergence.

We also notice of a risk in QKNorm, *i.e.*, there is a factor \sqrt{d} in both $P_{ij}^{(2)}$ and $\frac{\partial P_{ij}^{(2)}}{\partial \mathbf{x}_i}$. Therefore, with the increase of model size, there is still some potential risk of instability in training.

3.3.2 *StableAtten*

Based on the above analysis, we present a stabilized attention mechanism, called *StableAtten* as,

$$\mathbf{A} = \text{softmax}(\tau \mathbf{P}^{(3)}), \quad \text{where } P_{ij}^{(3)} = \left(\frac{SN(\mathbf{W}_q \mathbf{x}_i)^\top SN(\mathbf{W}_k \mathbf{x}_j)}{d_1^{2\alpha}} \right), \quad (3)$$

where τ is a temperature coefficient, and $SN(\cdot)$ is a *StableNorm* parameterized by α . Since that τ is likely related to the sequence length N , other than the head dimension d_1 , thus we set $\tau = 1.618 \cdot \log_2(N)$. When the input sequence length is 512, $\tau = 14.562$. We have that

$$P_{ij}^{(3)} = \left(\frac{\mathbf{W}_q \mathbf{x}_i}{\sqrt{\|\mathbf{W}_q \mathbf{x}_i\|_2^2 + \epsilon}} \right)^\top \text{diag}(\gamma_q) \text{diag}(\gamma_k) \frac{\mathbf{W}_k \mathbf{x}_j}{\sqrt{\|\mathbf{W}_k \mathbf{x}_j\|_2^2 + \epsilon}}.$$

The advantage of our stabilized form is that the logit is no longer directly related to the hidden dimension d . Although d is large in big model, the range of the logit in our *StableAtten* will not increase with d . Each *StableNorm* in query and key normalization will bring in a d^α , by dividing $d^{2\alpha}$, we can remove the influence of d .

3.3.3 EVALUATION FOR *StableAtten*

Here we will evaluate our *StableAtten*, as defined in Eq. 3, and compare it with the original dot-product attention. We evaluate on ViT-Large. The results are shown in Figure 4. We can see that from Figure 4:

- From the left panel of Figure 4, we see that increasing the number of heads in attention will improve the stability of the model training;

486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539

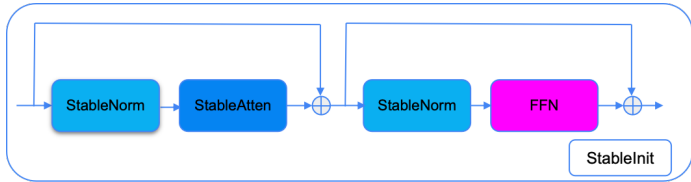


FIGURE 5: One block of our Stable Transformer. Our Stable Transformer uses StableNorm to replace LayerNorm or RMSNorm, use StableAtten to replace the original dot-product attention, and use StableInit to initialize the weights.

- From the right panel of Figure 4, we see that ViT-Large with *StableAtten* are more stable than their corresponding counterparts.

3.4 STABLE-TRANSFORMER

Following the architecture of Transformer, we can build a stabilized Transformer. As shown in Figure 5, we use *StableNorm* to replace LayerNorm or RMSNorm, use *StableAtten* to replace the original dot-product attention, and keep the same FFN module. Moreover, we use *StableInit* to initialize the weights.

In practice, our *Stable-Transformer* can lead to two different architectures, *i.e.*, a pure encoder architecture, *i.e.*, Vision Transformer, and a pure decoder architecture, *i.e.*, GPT. Accordingly, we name them as *StableViT* and *StableGPT*, respectively.

To be more specific, for ViT, we build four variants of *StableViT* in correspondence with ViT as detailed in Table 1 in the appendix: *StableViT-Large*, *StableViT-Huge*, *StableViT-giant*, and *StableViT-200* (which scales of parameters range from 307M to 1.44B); for GPT, we build three variants of *StableGPT* in correspondence with GPT as detailed in Table 2 in the appendix: *StableGPT-Small*, *StableGPT-Medium* and *StableGPT-Large* (which scales of parameters range from 124M to 774M).

3.4.1 EVALUATION FOR *Stable-Transformer*

To verify the effectiveness of our stabilized architecture, we evaluate different variants of *StableViT* and *Stable-GPT*. The experimental configurations of *StableViT* and *StableGPT* are shown in Table 3 and the experimental results are shown in Figure 1. We find that: our *StableGPT* can be trained smoothly, achieving a better validation loss (*i.e.*, 2.827 versus 2.848) compared to the original GPT2; our *StableViT* yields a better recognition accuracy (*i.e.*, 82.4% versus 81.3%) compared to the original ViT.

Moreover, we note that our *StableViT* and *StableGPT* can also tolerate larger learning rate. More empirical results and details are provided in Appendix I.

4 CONCLUSION

We have presented a theoretical analysis for the initialization, normalization and attention module of Transformer from the perspective of training stability. Specifically, we derived an upper bound and a lower bound for the expectations of the maximum and the minimum of the singular values of weight matrix obtained from Xavier initialization, found the reason why increasing hidden dimension can make the normalization layer likely leading to training instability from the Lipschits constant of the Jacobian matrix of the normalization layer, and also pointed out the theoretical mechanism why the hidden dimension can bring instability issue to affect self-attention module. Accordingly, we proposed three stabilized counterpart designs, *i.e.*, *StableInit*, *StableNorm* and *StableAtten*, and by putting them together, we also proposed a *Stable-Transformer*. We compiled our stabilized components and *Stable-Transformer* with GPT and ViT, and demonstrated that our stabilized methods can improve the training stability, leading improved performance. We hope that our work can benefit the deployment of larger deep models, especially the large language models, in varied application scenarios.

540 ETHICS STATEMENT

541

542 In this paper, we aim to provide a stabilized transformer. Our work does not involve any hu-
 543 man subjects, and we have carefully ensured that it poses no potential risks or harms. Addition-
 544 ally, there are no conflicts of interest, sponsorship concerns, or issues related to discrimination,
 545 bias, or fairness associated with this study. We have taken steps to address privacy and security
 546 concerns, and all data used comply with legal and ethical standards. Our work fully adheres to
 547 research integrity principles, and no ethical concerns have arisen during the course of this study.

548

549 REPRODUCIBILITY STATEMENT

550

551 To ensure the reproducibility of our work, we provide all the details to reproduce the experi-
 552 ments. Theoretical proofs of the claims made in this paper, and detailed experimental settings
 553 and configurations are provided in Appendices.

554

555 REFERENCES

556

557 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman,
 558 Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical
 559 report. *arXiv preprint arXiv:2303.08774*, 2023.

560

561 Devansh Arpit, Víctor Campos, and Yoshua Bengio. How to initialize your network? robust ini-
 562 tialization for weightnorm & resnets. *Advances in Neural Information Processing Systems*, 32,
 563 2019.

564

565 Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint*
 566 *arXiv:1607.06450*, 2016.

567

568 Thomas Bachlechner, Bodhisattwa Prasad Majumder, Henry Mao, Gary Cottrell, and Julian
 569 McAuley. Rezero is all you need: Fast convergence at large depth. In *Uncertainty in Artificial*
 570 *Intelligence*, pp. 1352–1361. PMLR, 2021.

571

572 Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly
 573 learning to align and translate. *arXiv preprint arXiv:1409.0473*, 2014.

574

575 Yoshua Bengio, Yann Lecun, and Geoffrey Hinton. Deep learning for ai. *Communications of the*
 576 *ACM*, 64(7):58–65, 2021.

577

578 Jeremy Bernstein, Arash Vahdat, Yisong Yue, and Ming-Yu Liu. On the distance between two neu-
 579 ral networks and the stability of learning. *Advances in Neural Information Processing Systems*,
 33:21370–21381, 2020.

580

581 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,
 582 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are
 583 few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.

584

585 Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam
 586 Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm:
 587 Scaling language modeling with pathways. *Journal of Machine Learning Research*, 24(240):
 1–113, 2023.

588

589 Aaron Defazio, Harsh Mehta, Konstantin Mishchenko, Ahmed Khaled, Ashok Cutkosky, et al. The
 590 road less scheduled. *arXiv preprint arXiv:2405.15682*, 2024.

591

592 Mostafa Dehghani, Josip Djolonga, Basil Mustafa, Piotr Padlewski, Jonathan Heek, Justin Gilmer,
 593 Andreas Peter Steiner, Mathilde Caron, Robert Geirhos, Ibrahim Alabdulmohsin, et al. Scaling
 vision transformers to 22 billion parameters. In *International Conference on Machine Learning*,
 pp. 7480–7512. PMLR, 2023.

- 594 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas
595 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An
596 image is worth 16x16 words: Transformers for image recognition at scale. In *International
597 Conference on Learning Representations*, 2020.
- 598 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha
599 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models.
600 *arXiv preprint arXiv:2407.21783*, 2024.
- 601 Alan Edelman and N Raj Rao. Random matrix theory. *Acta numerica*, 14:233–297, 2005.
- 602 Mahyar Fazlyab, Alexander Robey, Hamed Hassani, Manfred Morari, and George Pappas. Effi-
603 cient and accurate estimation of lipschitz constants for deep neural networks. *Advances in
604 Neural Information Processing Systems*, 32, 2019.
- 605 Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neu-
606 ral networks. In *Proceedings of the thirteenth international conference on artificial intelligence
607 and statistics*, pp. 249–256. JMLR Workshop and Conference Proceedings, 2010.
- 608 Ian Goodfellow, Yoshua Bengio, Aaron Courville, and Yoshua Bengio. *Deep learning*, volume 1.
609 MIT Press, 2016.
- 610 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing
611 human-level performance on imagenet classification. In *Proceedings of the IEEE international
612 conference on computer vision*, pp. 1026–1034, 2015.
- 613 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recog-
614 nition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp.
615 770–778, 2016.
- 616 Alex Henry, Prudhvi Raj Dachapally, Shubham Shantaram Pawar, and Yuxuan Chen. Query-key
617 normalization for transformers. In *Findings of the Association for Computational Linguistics:
618 EMNLP 2020*, pp. 4246–4253, 2020.
- 619 Xiao Shi Huang, Felipe Perez, Jimmy Ba, and Maksims Volkovs. Improving transformer opti-
620 mization through better initialization. In *International Conference on Machine Learning*, pp.
621 4475–4483. PMLR, 2020.
- 622 Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by
623 reducing internal covariate shift. In *International conference on machine learning*, pp. 448–456.
624 PMLR, 2015.
- 625 Andrej Karpathy. NanoGPT. <https://github.com/karpathy/nanoGPT>, 2022.
- 626 Hyunjik Kim, George Papamakarios, and Andriy Mnih. The lipschitz constant of self-attention.
627 In *International Conference on Machine Learning*, pp. 5562–5571. PMLR, 2021.
- 628 Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint
629 arXiv:1412.6980*, 2014.
- 630 Tim Large, Yang Liu, Minyoung Huh, Hyojin Bahng, Phillip Isola, and Jeremy Bernstein. Scalable
631 optimization in the modular norm. *arXiv preprint arXiv:2405.14813*, 2024.
- 632 Yann LeCun, Bernhard Boser, John S Denker, Donnie Henderson, Richard E Howard, Wayne Hub-
633 bard, and Lawrence D Jackel. Backpropagation applied to handwritten zip code recognition.
634 *Neural computation*, 1(4):541–551, 1989.
- 635 Tim Large, Yang Liu, Minyoung Huh, Hyojin Bahng, Phillip Isola, and Jeremy Bernstein. Scalable
636 optimization in the modular norm. *arXiv preprint arXiv:2405.14813*, 2024.
- 637 Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied
638 to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- 639 Yann LeCun, Léon Bottou, Genevieve B Orr, and Klaus-Robert Müller. Efficient backprop. In
640 *Neural networks: Tricks of the trade*, pp. 9–50. Springer, 2002.
- 641 Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *nature*, 521(7553):436–444,
642 2015.

- 648 Hong Liu, Zhiyuan Li, David Hall, Percy Liang, and Tengyu Ma. Sophia: A scalable stochastic
649 second-order optimizer for language model pre-training. *arXiv preprint arXiv:2305.14342*,
650 2023.
- 651 Liyuan Liu, Xiaodong Liu, Jianfeng Gao, Weizhu Chen, and Jiawei Han. Understanding the diffi-
652 culty of training transformers. In *Proceedings of the 2020 Conference on Empirical Methods in*
653 *Natural Language Processing (EMNLP)*, pp. 5747–5763, 2020.
- 654 Ze Liu, Han Hu, Yutong Lin, Zhuliang Yao, Zhenda Xie, Yixuan Wei, Jia Ning, Yue Cao, Zheng
655 Zhang, Li Dong, et al. Swin transformer v2: Scaling up capacity and resolution. In *Proceed-*
656 *ings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 12009–12019,
657 2022.
- 658 Ilya Loshchilov and Frank Hutter. Sgdr: Stochastic gradient descent with warm restarts. *arXiv*
659 *preprint arXiv:1608.03983*, 2016.
- 660 Ilya Loshchilov and Frank Hutter. Fixing weight decay regularization in adam. In *International*
661 *Conference on Learning Representations*, 2019.
- 662 Song Mei. Statistics 210b. <https://pillowmath.github.io/>, Spring 2022. Course Lec-
663 ture Notes.
- 664 Takeru Miyato, Toshiki Kataoka, Masanori Koyama, and Yuichi Yoshida. Spectral normalization
665 for generative adversarial networks. In *International Conference on Learning Representations*,
666 2018.
- 667 Vinod Nair and Geoffrey E Hinton. Rectified linear units improve restricted boltzmann machines.
668 In *Proceedings of the 27th international conference on machine learning (ICML-10)*, pp. 807–
669 814, 2010.
- 670 Yurii Nesterov. A method for unconstrained convex minimization problem with the rate of con-
671 vergence $O(1/k^2)$. In *Doklady an ussr*, volume 269, pp. 543–547, 1983.
- 672 Yurri Nesterov. Introductory lectures on convex programming volume i: Basic course. 1998.
- 673 Toan Q Nguyen and Julian Salazar. Transformers without tears: Improving the normalization of
674 self-attention. *arXiv preprint arXiv:1910.05895*, 2019.
- 675 Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor
676 Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-
677 performance deep learning library. *Advances in neural information processing systems*, 32,
678 2019.
- 679 William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of*
680 *the IEEE/CVF International Conference on Computer Vision*, pp. 4195–4205, 2023.
- 681 Xianbiao Qi, Jianan Wang, Yihao Chen, Yukai Shi, and Lei Zhang. Lipsformer: Introducing lips-
682 chitz continuity to vision transformers. In *The Eleventh International Conference on Learning*
683 *Representations*, 2023a.
- 684 Xianbiao Qi, Jianan Wang, and Lei Zhang. Understanding optimization of deep learning via ja-
685 cobian matrix and lipschitz constant. *arXiv preprint arXiv:2306.09338*, 2023b.
- 686 Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. Improving language un-
687 derstanding by generative pre-training. 2018.
- 688 Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language
689 models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- 690 Nikhila Ravi, Valentin Gabeur, Yuan-Ting Hu, Ronghang Hu, Chaitanya Ryali, Tengyu Ma,
691 Haitham Khedr, Roman Rädle, Chloe Rolland, Laura Gustafson, et al. Sam 2: Segment any-
692 thing in images and videos. *arXiv preprint arXiv:2408.00714*, 2024.
- 693 David E Rumelhart, Geoffrey E Hinton, and Ronald J Williams. Learning representations by back-
694 propagating errors. *nature*, 323(6088):533–536, 1986.

- 702 Tim Salimans and Durk P Kingma. Weight normalization: A simple reparameterization to accel-
703 erate training of deep neural networks. *Advances in neural information processing systems*, 29,
704 2016.
- 705 Andrew M Saxe, James L McClelland, and Surya Ganguli. Exact solutions to the nonlinear dy-
706 namics of learning in deep linear neural networks. *arXiv preprint arXiv:1312.6120*, 2013.
- 707 Ilya Sutskever, James Martens, George Dahl, and Geoffrey Hinton. On the importance of initial-
708 ization and momentum in deep learning. In *International conference on machine learning*, pp.
709 1139–1147. PMLR, 2013.
- 710 Terence Tao. *Topics in random matrix theory*, volume 132. American Mathematical Soc., 2012.
- 711 Qwen Team. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023.
- 712 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Niko-
713 lay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open founda-
714 tion and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- 715 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
716 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural informa-
717 tion processing systems*, 30, 2017.
- 718 Roman Vershynin. Introduction to the non-asymptotic analysis of random matrices. *arXiv
719 preprint arXiv:1011.3027*, 2010.
- 720 Roman Vershynin. *High-dimensional probability: An introduction with applications in data sci-
721 ence*, volume 47. Cambridge university press, 2018.
- 722 Hongyu Wang, Shuming Ma, Li Dong, Shaohan Huang, Dongdong Zhang, and Furu Wei. Deep-
723 net: Scaling transformers to 1,000 layers. *arXiv preprint arXiv:2203.00555*, 2022.
- 724 Qiang Wang, Bei Li, Tong Xiao, Jingbo Zhu, Changliang Li, Derek F Wong, and Lidia S Chao.
725 Learning deep transformer models for machine translation. *arXiv preprint arXiv:1906.01787*,
726 2019.
- 727 Ross Wightman. Pytorch image models. [https://github.com/rwightman/
728 pytorch-image-models](https://github.com/rwightman/pytorch-image-models), 2019.
- 729 Eugene P Wigner. Characteristic vectors of bordered matrices with infinite dimensions. *Annals
730 of Mathematics*, 62(3):548–564, 1955.
- 731 Mitchell Wortsman, Peter J Liu, Lechao Xiao, Katie E Everett, Alexander A Alemi, Ben Adlam,
732 John D Co-Reyes, Izzeddin Gur, Abhishek Kumar, Roman Novak, et al. Small-scale proxies
733 for large-scale transformer training instabilities. In *The Twelfth International Conference on
734 Learning Representations*, 2024.
- 735 Xingyu Xie, Pan Zhou, Huan Li, Zhouchen Lin, and Shuicheng Yan. Adan: Adaptive nesterov mo-
736 mentum algorithm for faster optimizing deep models. *IEEE Transactions on Pattern Analysis
737 and Machine Intelligence*, 2024.
- 738 Ruibin Xiong, Yunchang Yang, Di He, Kai Zheng, Shuxin Zheng, Chen Xing, Huishuai Zhang,
739 Yanyan Lan, Liwei Wang, and Tiejian Liu. On layer normalization in the transformer architec-
740 ture. In *International Conference on Machine Learning*, pp. 10524–10533. PMLR, 2020.
- 741 Greg Yang. Wide feedforward or recurrent neural networks of any architecture are gaussian pro-
742 cesses. *Advances in Neural Information Processing Systems*, 32, 2019.
- 743 Greg Yang and Edward J Hu. Tensor programs iv: Feature learning in infinite-width neural net-
744 works. In *International Conference on Machine Learning*, pp. 11727–11737. PMLR, 2021.
- 745 Greg Yang, Edward J Hu, Igor Babuschkin, Szymon Sidor, Xiaodong Liu, David Farhi, Nick Ry-
746 der, Jakub Pachocki, Weizhu Chen, and Jianfeng Gao. Tensor programs v: Tuning large neural
747 networks via zero-shot hyperparameter transfer. *arXiv preprint arXiv:2203.03466*, 2022.
- 748
749
750
751
752
753
754
755

756 Biao Zhang and Rico Sennrich. Root mean square layer normalization. *Advances in Neural Infor-*
757 *mation Processing Systems*, 32, 2019.
758
759 Hongyi Zhang, Yann N Dauphin, and Tengyu Ma. Fixup initialization: Residual learning without
760 normalization. *arXiv preprint arXiv:1901.09321*, 2019.
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809

A NOTATIONS

We primarily follow the notations used in the renowned deep learning book (Goodfellow et al., 2016). We use **bold symbol** to denote a column vector or a matrix, and use non-bold symbol to denote scalar. For instance $\mathbf{y} = \mathbf{W}\mathbf{x}$, where \mathbf{x} and \mathbf{y} are two column vectors and \mathbf{W} is a projection matrix. We use denominator layout³, the Jacobian matrix of \mathbf{y} with respect to \mathbf{x} is $\frac{\partial \mathbf{W}\mathbf{x}}{\partial \mathbf{x}} = \mathbf{W}^\top$, and we have $\frac{\partial \mathbf{x}^\top \mathbf{W}\mathbf{x}}{\partial \mathbf{x}} = (\mathbf{W} + \mathbf{W}^\top)\mathbf{x}$. Using the denominator layout, for a chain function $o = f(g(h(\mathbf{x})))$, where $\mathbf{y} = h(\mathbf{x})$, $\mathbf{z} = g(\mathbf{y})$, and $o = f(\mathbf{z})$. we have the Jacobian matrix of o with respect to \mathbf{x} as $\frac{\partial o}{\partial \mathbf{x}} = \frac{\partial \mathbf{y}}{\partial \mathbf{x}} \frac{\partial \mathbf{z}}{\partial \mathbf{y}} \frac{\partial o}{\partial \mathbf{z}}$.

B PROOF OF THEOREM 1

Our proof of Theorem 1 is based on references (Vershynin, 2018; 2010; Mei, Spring 2022). Let us first clarify our problem, we have $\mathbf{W} \in \mathcal{R}^{m \times n}$, where $W_{i,j} \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1)$, we need to prove $\sqrt{m} - \sqrt{n} \leq \mathbb{E}[\sigma_{\min}(\mathbf{W})] \leq \mathbb{E}[\sigma_{\max}(\mathbf{W})] \leq \sqrt{m} + \sqrt{n}$. To prove Theorem 1, we need to prove two parts, i.e., $\mathbb{E}[\sigma_{\max}(\mathbf{W})] \leq \sqrt{m} + \sqrt{n}$ and $\sqrt{m} - \sqrt{n} \leq \mathbb{E}[\sigma_{\min}(\mathbf{W})]$. To prove the first part, we need to first introduce Sudakov-Fernique inequality.

Theorem 4 (Sudakov-Fernique inequality.)

Let $(A_{s,t})_{s \in S, t \in T}$ and $(B_{s,t})_{s \in S, t \in T}$ be two zero mean Gaussian processes. Assume that for all $s_1, s_2 \in S$ and $t_1, t_2 \in T$, we have

$$\mathbb{E}[(A_{t_1, s_1} - A_{t_2, s_2})^2] \leq \mathbb{E}[(B_{t_1, s_1} - B_{t_2, s_2})^2].$$

Then we have

$$\mathbb{E} \left[\sup_{s \in S, t \in T} A_{s,t} \right] \leq \mathbb{E} \left[\sup_{s \in S, t \in T} B_{s,t} \right].$$

Here, we do not provide the proof of Theorem 4. You can find the proof in (Vershynin, 2018). The Sudakov-Fernique inequality will be used in our following proof of $\mathbb{E}[\sigma_{\max}(\mathbf{W})] \leq \sqrt{m} + \sqrt{n}$.

Let us define $A_{\mathbf{u}, \mathbf{v}} = \langle \mathbf{W}\mathbf{u}, \mathbf{v} \rangle = \mathbf{v}^\top \mathbf{W}\mathbf{u}$ for $\mathbf{u} \in S^{n-1}$ and $\mathbf{v} \in S^{m-1}$, where $W_{i,j} \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1)$. We define

$$B_{\mathbf{u}, \mathbf{v}} = \langle \mathbf{u}, \mathbf{g} \rangle + \langle \mathbf{v}, \mathbf{h} \rangle = \sum_{i=1}^n u_i g_i + \sum_{j=1}^m v_j g_j, \quad g_i \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1), h_j \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1).$$

For any $(\mathbf{u}, \mathbf{v}), (\mathbf{q}, \mathbf{z}) \in (S^{n-1} \times S^{m-1})$, let us consider that:

$$\begin{aligned} \mathbb{E} \left[(A_{\mathbf{u}, \mathbf{v}} - A_{\mathbf{q}, \mathbf{z}})^2 \right] &= \mathbb{E} \left[(\langle \mathbf{W}\mathbf{u}, \mathbf{v} \rangle - \langle \mathbf{W}\mathbf{q}, \mathbf{z} \rangle)^2 \right] \\ &= \mathbb{E} \left[\left(\sum_{i,j} W_{i,j} (u_j v_i - q_j z_i) \right)^2 \right] \\ &= \sum_{i,j} (u_j v_i - q_j z_i)^2 \quad (\text{by independence, } W_{i,j} \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1)) \\ &= \|\mathbf{u}\mathbf{v}^\top - \mathbf{q}\mathbf{z}^\top\|_F^2 \\ &\leq \|\mathbf{u} - \mathbf{q}\|_2^2 + \|\mathbf{v} - \mathbf{z}\|_2^2. \quad (\text{see the following proof.}) \end{aligned}$$

We need to prove the following inequality:

$$\|\mathbf{u}\mathbf{v}^\top - \mathbf{q}\mathbf{z}^\top\|_F^2 \leq \|\mathbf{u} - \mathbf{q}\|_2^2 + \|\mathbf{v} - \mathbf{z}\|_2^2,$$

where $\mathbf{u}, \mathbf{q} \in S^{n-1}$ and $\mathbf{v}, \mathbf{z} \in S^{m-1}$, and S^{n-1} denotes the unit sphere in \mathbb{R}^n .

³https://en.wikipedia.org/wiki/Matrix_calculus

864 *Proof.* First, recall that $\|\mathbf{W}\|_F^2 = \text{tr}[\mathbf{W}^\top \mathbf{W}]$, we have:

$$\begin{aligned}
865 & \\
866 & \|\mathbf{u}\mathbf{v}^\top - \mathbf{q}\mathbf{z}^\top\|_F^2 = \text{tr}[(\mathbf{u}\mathbf{v}^\top - \mathbf{q}\mathbf{z}^\top)^\top (\mathbf{u}\mathbf{v}^\top - \mathbf{q}\mathbf{z}^\top)]. \\
867 & = \text{tr}[(\mathbf{v}\mathbf{u}^\top - \mathbf{z}\mathbf{q}^\top)(\mathbf{u}\mathbf{v}^\top - \mathbf{q}\mathbf{z}^\top)]. \\
868 & = \text{tr}[\mathbf{v}(\mathbf{u}^\top \mathbf{u})\mathbf{v}^\top - \mathbf{v}(\mathbf{u}^\top \mathbf{q})\mathbf{z}^\top - \mathbf{z}(\mathbf{q}^\top \mathbf{u})\mathbf{v}^\top + \mathbf{z}(\mathbf{q}^\top \mathbf{q})\mathbf{z}^\top]. \\
869 & = \text{tr}[\mathbf{v}\mathbf{v}^\top - \mathbf{v}\mathbf{z}^\top \mathbf{q}^\top \mathbf{u} - \mathbf{z}\mathbf{u}^\top \mathbf{q}^\top \mathbf{v} + \mathbf{z}\mathbf{z}^\top]. \\
870 & = \text{tr}(\mathbf{v}\mathbf{v}^\top) + \text{tr}(\mathbf{z}\mathbf{z}^\top) - \text{tr}(\mathbf{v}\mathbf{z}^\top \mathbf{q}^\top \mathbf{u}) - \text{tr}(\mathbf{z}\mathbf{u}^\top \mathbf{q}^\top \mathbf{v}). \\
871 & \\
872 &
\end{aligned}$$

873 From the definition of the trace, we have:

$$874 \text{tr}(\mathbf{v}\mathbf{v}^\top) = \|\mathbf{v}\|_2^2, \quad \text{tr}(\mathbf{z}\mathbf{z}^\top) = \|\mathbf{z}\|_2^2, \quad \text{tr}(\mathbf{v}\mathbf{z}^\top \mathbf{q}^\top \mathbf{u}) = (\mathbf{v}^\top \mathbf{z})(\mathbf{q}^\top \mathbf{u}), \quad \text{tr}(\mathbf{z}\mathbf{u}^\top \mathbf{q}^\top \mathbf{v}) = (\mathbf{z}^\top \mathbf{v})(\mathbf{u}^\top \mathbf{q}).$$

875 Thus, we have:

$$876 \|\mathbf{u}\mathbf{v}^\top - \mathbf{q}\mathbf{z}^\top\|_F^2 = \|\mathbf{v}\|_2^2 \|\mathbf{u}\|_2^2 + \|\mathbf{z}\|_2^2 \|\mathbf{q}\|_2^2 - 2(\mathbf{v}^\top \mathbf{z})(\mathbf{u}^\top \mathbf{q}).$$

877 Since $\mathbf{u}, \mathbf{q} \in S^{n-1}$ and $\mathbf{v}, \mathbf{z} \in S^{m-1}$, we have $\|\mathbf{u}\|_2 = \|\mathbf{q}\|_2 = 1$ and $\|\mathbf{v}\|_2 = \|\mathbf{z}\|_2 = 1$, simplifying to:

$$\begin{aligned}
878 & \|\mathbf{v}\|_2^2 \|\mathbf{u}\|_2^2 + \|\mathbf{z}\|_2^2 \|\mathbf{q}\|_2^2 - 2(\mathbf{v}^\top \mathbf{z})(\mathbf{u}^\top \mathbf{q}) = 1 + 1 - 2(\mathbf{v}^\top \mathbf{z})(\mathbf{u}^\top \mathbf{q}). \\
879 & = 2 - 2(\mathbf{v}^\top \mathbf{z})(\mathbf{u}^\top \mathbf{q}). \\
880 &
\end{aligned}$$

881 Now, consider the right-hand side:

$$\begin{aligned}
882 & \|\mathbf{u} - \mathbf{q}\|_2^2 + \|\mathbf{v} - \mathbf{z}\|_2^2 = (\|\mathbf{u}\|_2^2 - 2\mathbf{u}^\top \mathbf{q} + \|\mathbf{q}\|_2^2) + (\|\mathbf{v}\|_2^2 - 2\mathbf{v}^\top \mathbf{z} + \|\mathbf{z}\|_2^2). \\
883 & = 1 - 2\mathbf{u}^\top \mathbf{q} + 1 + 1 - 2\mathbf{v}^\top \mathbf{z} + 1. \\
884 & = 2 - 2(\mathbf{u}^\top \mathbf{q}) + 2 - 2(\mathbf{v}^\top \mathbf{z}). \\
885 & = 4 - 2(\mathbf{u}^\top \mathbf{q}) - 2(\mathbf{v}^\top \mathbf{z}). \\
886 &
\end{aligned}$$

887 let us assume $c = \mathbf{u}^\top \mathbf{q}$ and $d = \mathbf{v}^\top \mathbf{z}$. Since $\mathbf{u}, \mathbf{q} \in S^{n-1}$ and $\mathbf{v}, \mathbf{z} \in S^{m-1}$, we know $c \leq 1$, $d \leq 1$, and

$$888 (4 - 2(\mathbf{u}^\top \mathbf{q}) - 2(\mathbf{v}^\top \mathbf{z})) - (2 - 2(\mathbf{u}^\top \mathbf{q})(\mathbf{v}^\top \mathbf{z})) = 2(1 - c)(1 - d) \geq 0.$$

889 Similarly, we have,

$$\begin{aligned}
890 & \mathbb{E}[(B_{\mathbf{u},\mathbf{v}} - B_{\mathbf{q},\mathbf{z}})^2] = \mathbb{E}[(\langle \mathbf{g}, \mathbf{u} - \mathbf{q} \rangle + \langle \mathbf{h}, \mathbf{v} - \mathbf{z} \rangle)^2] \\
891 & = \mathbb{E}[\langle \mathbf{g}, \mathbf{u} - \mathbf{q} \rangle^2] + \mathbb{E}[\langle \mathbf{h}, \mathbf{v} - \mathbf{z} \rangle^2] \quad (\text{by independence, mean 0}) \\
892 & = \|\mathbf{u} - \mathbf{q}\|_2^2 + \|\mathbf{v} - \mathbf{z}\|_2^2. \quad (\text{since } \mathbf{g}, \mathbf{h} \text{ are standard normal}). \\
893 &
\end{aligned}$$

894 Thus, we have

$$895 \mathbb{E}[A_{\mathbf{u},\mathbf{v}} - A_{\mathbf{q},\mathbf{z}}]^2 \leq \mathbb{E}[(B_{\mathbf{u},\mathbf{v}} - B_{\mathbf{q},\mathbf{z}})^2].$$

896 Now, applying the Sudakov-Fernique inequality, we have

$$\begin{aligned}
897 & \mathbb{E} \left[\sup_{(\mathbf{u}, \mathbf{v}) \in S^{n-1} \times S^{m-1}} \langle \mathbf{W}\mathbf{u}, \mathbf{v} \rangle \right] \leq \mathbb{E} \left[\sup_{(\mathbf{u}, \mathbf{v}) \in S^{n-1} \times S^{m-1}} (\langle \mathbf{u}, \mathbf{g} \rangle + \langle \mathbf{v}, \mathbf{h} \rangle) \right] \\
898 & = \mathbb{E} \left[\sup_{(\mathbf{u}, \mathbf{v}) \in S^{n-1} \times S^{m-1}} \langle \mathbf{u}, \mathbf{g} \rangle \right] + \mathbb{E} \left[\sup_{(\mathbf{u}, \mathbf{v}) \in S^{n-1} \times S^{m-1}} \langle \mathbf{v}, \mathbf{h} \rangle \right] \\
899 & = \mathbb{E}[\|\mathbf{g}\|_2] + \mathbb{E}[\|\mathbf{h}\|_2] \\
900 & \leq \mathbb{E}[\|\mathbf{g}\|_2^2]^{1/2} + \mathbb{E}[\|\mathbf{h}\|_2^2]^{1/2} \\
901 & = \sqrt{n} + \sqrt{m}. \\
902 &
\end{aligned}$$

903 This completes the proof of the first part of Theorem 1.

918

□

919

920

To prove the second part, we need to introduce Gordon's Inequality.

921

Theorem 5 (Gordon's Inequality.)

922

923

Let $(A_{s,t})_{s \in S, t \in T}$ and $(B_{s,t})_{s \in S, t \in T}$ be two Gaussian processes with $\mathbb{E}[A_{s,t}] = \mathbb{E}[B_{s,t}] = 0$, and suppose that

925

926

927

$$\begin{cases} \mathbb{E}[(A_{s_1, t_1} - A_{s_2, t_2})^2] \geq \mathbb{E}[(B_{s_1, t_1} - B_{s_2, t_2})^2] & \forall t_1, t_2 \in T, s \in S, \\ \mathbb{E}[(A_{s_1, t_1} - A_{s_2, t_2})^2] \leq \mathbb{E}[(B_{s_1, t_1} - B_{s_2, t_2})^2] & \forall s_1 \neq s_2 \in S, t_1, t_2 \in T. \end{cases}$$

928

Then

929

930

$$\mathbb{E} \left[\sup_{s \in S} \inf_{t \in T} A_{s,t} \right] \leq \mathbb{E} \left[\sup_{s \in S} \inf_{t \in T} B_{s,t} \right].$$

931

932

933

Same as above, we do not provide the proof of Gordon's inequality. The proof can be found in (Vershynin, 2018). We will directly use it to help our proof of Theorem 1.

934

935

936

Proof. Let $B_{u,v} = \langle \mathbf{g}, \mathbf{u} \rangle + \langle \mathbf{h}, \mathbf{v} \rangle$. Check that $A_{u,v}$ and $B_{u,v}$ satisfy the conditions in the theorem. Then we have

937

938

939

940

941

942

943

944

945

946

947

948

949

950

951

$$\begin{aligned} -\mathbb{E}[\sigma_{\min}(\mathbf{W})] &= \mathbb{E} \left[\sup_{\mathbf{v} \in S^{n-1}} -\|\mathbf{W}\mathbf{v}\|_2 \right] \\ &= \mathbb{E} \left[\sup_{\mathbf{v} \in S^{n-1}} \inf_{\mathbf{u} \in S^{m-1}} \langle \mathbf{u}, -\mathbf{W}\mathbf{v} \rangle \right] \\ &\leq \mathbb{E} \left[\sup_{\mathbf{v} \in S^{n-1}} \inf_{\mathbf{u} \in S^{m-1}} \langle \mathbf{g}, \mathbf{u} \rangle + \langle \mathbf{h}, \mathbf{v} \rangle \right] \\ &= \mathbb{E} \left[\sup_{\mathbf{v} \in S^{n-1}} \langle \mathbf{h}, \mathbf{v} \rangle \right] + \mathbb{E} \left[\inf_{\mathbf{u} \in S^{m-1}} \langle \mathbf{g}, \mathbf{u} \rangle \right] \quad (\text{since } \mathbf{g}, \mathbf{h} \text{ are standard normal.}) \\ &= \mathbb{E}[\|\mathbf{h}\|_2] - \mathbb{E}[\|\mathbf{g}\|_2] \\ &= \sqrt{n} - \sqrt{m}. \end{aligned}$$

952

Therefore, we have

953

$$\mathbb{E}[\sigma_{\min}(\mathbf{W})] \geq \sqrt{m} - \sqrt{n}.$$

954

This completes the proof of the second part of Theorem 1. □

955

956

957

958

C PROOF OF $\sigma_{\max} \left(\mathbf{I} - \frac{\mathbf{y}\mathbf{y}^T}{\|\mathbf{y}\|_2^2 + \epsilon} \right) \leq 1$

959

960

961

962

963

Proof. Let $\mathbf{M} = \left(\mathbf{I} - \frac{\mathbf{y}\mathbf{y}^T}{\|\mathbf{y}\|_2^2 + \epsilon} \right)$, to prove that the maximum singular value of the matrix \mathbf{M} is 1, we need to analyze the properties of this matrix.

964

965

966

Let $\mathbf{A} = \frac{\mathbf{y}\mathbf{y}^T}{\|\mathbf{y}\|_2^2 + \epsilon}$, the matrix \mathbf{A} is a rank-1 matrix with one non-zero eigenvalue. The non-zero eigenvalue is

967

968

$$\lambda_{\mathbf{A}} = \frac{\|\mathbf{y}\|_2^2}{\|\mathbf{y}\|_2^2 + \epsilon}$$

969

970

The eigenvalues of \mathbf{M} are $1 - \lambda_{\mathbf{A}}$ and 1 with multiplicity $n - 1$:

971

$$\lambda_{\mathbf{M}} = 1 - \frac{\|\mathbf{y}\|_2^2}{\|\mathbf{y}\|_2^2 + \epsilon} = \frac{\epsilon}{\|\mathbf{y}\|_2^2 + \epsilon}$$

All other eigenvalues are 1. The singular values of M are the absolute values of its eigenvalues:

$$\sigma_1(M) = 1, \quad (\text{with multiplicity } n-1)$$

$$\sigma_2(M) = \frac{\epsilon}{\|\mathbf{y}\|_2^2 + \epsilon}.$$

The maximum singular value of M is the largest eigenvalue in absolute value, which is 1. Thus, the maximum singular value of the matrix $\left(\mathbf{I} - \frac{\mathbf{y}\mathbf{y}^\top}{\|\mathbf{y}\|_2^2 + \epsilon}\right)$ is 1. \square

D QKNORM DERIVATIONS

Here, we list all the partial derivations for the QKNorm.

$$\begin{aligned} \frac{\partial P_{ij}^{(2)}}{\partial \mathbf{x}_i} &= \frac{\sqrt{d}}{\sqrt{\|\mathbf{W}_q \mathbf{x}_i\|_2^2 + \epsilon}} \mathbf{W}_q^\top \left(\mathbf{I} - \frac{\mathbf{W}_q \mathbf{x}_i (\mathbf{W}_q \mathbf{x}_i)^\top}{\|\mathbf{W}_q \mathbf{x}_i\|_2^2 + \epsilon} \right) \text{diag}(\gamma_q) \text{diag}(\gamma_k) \frac{\mathbf{W}_k \mathbf{x}_j}{\sqrt{\|\mathbf{W}_k \mathbf{x}_j\|_2^2 + \epsilon}}, \\ \frac{\partial P_{ij}^{(2)}}{\partial \mathbf{x}_j} &= \frac{\sqrt{d}}{\sqrt{\|\mathbf{W}_k \mathbf{x}_j\|_2^2 + \epsilon}} \mathbf{W}_k^\top \left(\mathbf{I} - \frac{\mathbf{W}_k \mathbf{x}_j (\mathbf{W}_k \mathbf{x}_j)^\top}{\|\mathbf{W}_k \mathbf{x}_j\|_2^2 + \epsilon} \right) \text{diag}(\gamma_k) \text{diag}(\gamma_q) \frac{\mathbf{W}_q \mathbf{x}_i}{\sqrt{\|\mathbf{W}_q \mathbf{x}_i\|_2^2 + \epsilon}}, \\ \frac{\partial P_{ij}^{(2)}}{\partial \mathbf{W}_q} &= \frac{\sqrt{d}}{\sqrt{\|\mathbf{W}_q \mathbf{x}_i\|_2^2 + \epsilon}} \left(\mathbf{I} - \frac{\mathbf{W}_q \mathbf{x}_i (\mathbf{W}_q \mathbf{x}_i)^\top}{\|\mathbf{W}_q \mathbf{x}_i\|_2^2 + \epsilon} \right) \text{diag}(\gamma_q) \text{diag}(\gamma_k) \frac{\mathbf{W}_k \mathbf{x}_j}{\sqrt{\|\mathbf{W}_k \mathbf{x}_j\|_2^2 + \epsilon}} \mathbf{x}_i^\top, \\ \frac{\partial P_{ij}^{(2)}}{\partial \mathbf{W}_k} &= \frac{\sqrt{d}}{\sqrt{\|\mathbf{W}_k \mathbf{x}_j\|_2^2 + \epsilon}} \left(\mathbf{I} - \frac{\mathbf{W}_k \mathbf{x}_j (\mathbf{W}_k \mathbf{x}_j)^\top}{\|\mathbf{W}_k \mathbf{x}_j\|_2^2 + \epsilon} \right) \text{diag}(\gamma_k) \text{diag}(\gamma_q) \frac{\mathbf{W}_q \mathbf{x}_i}{\sqrt{\|\mathbf{W}_q \mathbf{x}_i\|_2^2 + \epsilon}} \mathbf{x}_j^\top. \end{aligned}$$

When considering γ_q and γ_k are set to be 1, $\frac{\partial P_{ij}^{(2)}}{\partial \mathbf{x}_i}$ is only proportion to $\mathcal{O}(\|\mathbf{W}_q\|)$ and $\frac{\partial P_{ij}^{(2)}}{\partial \mathbf{W}_q} \leq \frac{\sqrt{d}}{\sqrt{\epsilon}}$.

E PROOF OF LEMMA 1

Proof. To prove $\mathbb{E}[\sigma_{\max}(\mathbf{W})] \leq 2$, it is equivalent to prove $\sqrt{\frac{2}{n_{in} + n_{out}}} (\sqrt{n_{in}} + \sqrt{n_{out}}) \leq 2$ for any n_{in} and n_{out} . Note that:

$$\left(\sqrt{\frac{2}{n_{in} + n_{out}}} (\sqrt{n_{in}} + \sqrt{n_{out}}) \right)^2 = \frac{2(\sqrt{n_{in}} + \sqrt{n_{out}})^2}{n_{in} + n_{out}} = \frac{2(n_{in} + n_{out}) + 4\sqrt{n_{in}n_{out}}}{n_{in} + n_{out}} = 2 + \frac{4\sqrt{n_{in}n_{out}}}{n_{in} + n_{out}} \leq 4.$$

Thus, we have $\left(\sqrt{\frac{2}{n_{in} + n_{out}}} (\sqrt{n_{in}} + \sqrt{n_{out}}) \right) \leq 2$. \square

F MODEL AND TRAINING CONFIGURATION

Model Configurations. We list some basic configurations of our *StableViT* and *StableGPT* in Table 1 and Table 2.

Training Configurations. We list the training configurations of our *StableGPT* and *StableViT* in Table 3. For *StableGPT*, we fully follow the experimental configurations of nanoGPT (Karpathy, 2022), all parameters are same as GPT2 (Radford et al., 2019). All experiments are conducted on

TABLE 1: Model configuration for *StableViT*. The *StableViT* is similar with the original ViT (Dosovitskiy et al., 2020).

Model Card	Params.	Blocks	Embed. dim.	MLP. dim.	Heads	Epochs	Peak LR
<i>StableViT-L-16</i>	307M	24	1024	4096	16	150 or 300	1e-3
<i>StableViT-H-14</i>	632M	32	1280	5120	16	150 or 300	1e-3
<i>StableViT-g-14</i>	1011M	40	1408	6144	16	150 or 300	1e-3
<i>StableViT-200</i>	1439M	200	768	3072	12	150 or 300	1e-3

TABLE 2: Model configuration for *StableGPT*. The *StableGPT* is similar with the original GPT2 (Radford et al., 2019). We do not include larger models as in nanoGPT (Karpathy, 2022) because training larger models will cost much more computational resource.

Model Card	Params.	Blocks	Embed. dim.	Heads	Train steps	Peak LR	Minimum LR
<i>StableGPT-S</i>	124M	12	768	12	600K	6e-4	6e-5
<i>StableGPT-M</i>	350M	24	1024	16	600K	3e-4	3e-5
<i>StableGPT-L</i>	774M	36	1280	20	600K	2.5e-4	2.5e-5

A800 GPU cluster. For instance, it takes around 3 days to train *StableGPT*-Small on a GPU server with 8 A800 GPUs. *StableGPT*-Medium will take around 7.5 days. Note that in the original ViT, we use 60 epochs’ learning rate warmup, but in our *StableViT*, we do not use warmup. We do not include some new optimizer (Liu et al., 2023) or learning schedule (Defazio et al., 2024) to further improve the performance of the models.

TABLE 3: Training configurations for *StableGPT* and *StableViT*.

	(a) Training configurations for <i>StableGPT</i> .		(b) Training configurations for <i>StableViT</i> .
training config	StableGPT-S/M/L	training config	StableViT-L/H/g/200 (224 ²)
weight init	<i>StableInit</i>	weight init	<i>StableInit</i>
optimizer	AdamW	optimizer	AdamW
baseline learning rate	0.0006	base learning rate	1e-3
weight decay	0.1	weight decay	0.1
optimizer momentum	$\beta_1, \beta_2 = 0.9, 0.95$	optimizer momentum	$\beta_1, \beta_2 = 0.9, 0.99$
warmup	2,000	batch size	1024
tokens seen each update	500,000	training epochs	300 or 150 or 60
max iters	600,000	learning rate schedule	cosine decay
batch size	480	warmup epochs	0
sequence length	1024	randaugment	(9, 0.5)
dropout	0.0	mixup	0.8
bfloat16	True	cutmix	1.0
gradient clipping	1.0	random erasing	0
		label smoothing	0.1
		stochastic depth	0.5/0.5
		gradient clip	None
		exp. mov. avg. (EMA)	no

G DEMONSTRATION CODE

To help the audience understand the details of the introduced modules, we list our demonstration codes.

CODE 1: Stable-Transformer Implementation Demonstration.

```

1086 1 import torch
1087 2 import torch.nn as nn
1088 3 import math
1089 4
1090 5 def StableInit(module: nn.Module, name: str = '') -> None:
1091 6     if isinstance(module, nn.Linear):
1092 7         n_in, n_out = module.weight.shape[0], module.weight.shape[1]
1093 8         init_std = 1.0/(math.sqrt(n_in)+math.sqrt(n_out))
1094 9         torch.nn.init.normal_(module.weight, mean=0.0, std=init_std)
1095 10        if module.bias is not None:
1096 11            nn.init.zeros_(module.bias)
1097 12
1098 13 class StableNorm(nn.Module):
1099 14     def __init__(self, ndim: int, alpha: float = 0.0, eps: float = 1e-8):
1100 15         super().__init__()
1101 16         self.alpha = alpha
1102 17         self.ndim = ndim
1103 18         self.eps = eps
1104 19         self.weight = nn.Parameter(torch.ones(ndim))
1105 20
1106 21     def forward(self, input):
1107 22         x_norm = torch.norm(input, dim=2, keepdim=True) + self.eps
1108 23         x = math.pow(self.ndim, self.alpha)*input/x_norm
1109 24         y = self.weight.unsqueeze(0).unsqueeze(0)*x
1110 25         return y
1111 26
1112 27 class StableAtten(nn.Module):
1113 28     def __init__(self, dim: int, num_heads: int = 8, qkv_bias: bool = False,
1114 29                 attn_drop: float = 0., proj_drop: float = 0.,
1115 30                 norm_layer: nn.Module = StableNorm,
1116 31                 temperature: float = 1.0, sequence_length: int=0) -> None:
1117 32         super().__init__()
1118 33         assert dim % num_heads == 0,
1119 34             self.num_heads = num_heads
1120 35             self.head_dim = dim // num_heads
1121 36             self.scale = self.head_dim ** -0.5
1122 37             self.qkv = nn.Linear(dim, dim * 3, bias=qkv_bias)
1123 38             self.q_norm = norm_layer(self.head_dim)
1124 39             self.k_norm = norm_layer(self.head_dim)
1125 40             norm_alpha = 2 * self.q_norm.alpha
1126 41             self.tau = 1.618*math.log(sequence_length,2)*temperature
1127 42             self.scale = self.head_dim**(-norm_alpha)*self.tau
1128 43             self.attn_drop = nn.Dropout(attn_drop)
1129 44             self.proj = nn.Linear(dim, dim)
1130 45             self.proj_drop = nn.Dropout(proj_drop)
1131 46
1132 47     def forward(self, x: torch.Tensor) -> torch.Tensor:
1133 48         B, N, C = x.shape
1134 49         qkv = self.qkv(x).reshape(B,N,3,self.num_heads,self.head_dim)
1135 50         qkv = qkv.permute(2,0,3,1,4)
1136 51         q, k, v = qkv.unbind(0)
1137 52         q, k = self.q_norm(q), self.k_norm(k)
1138 53         q = q * self.scale
1139 54         attn = q @ k.transpose(-2, -1)
1140 55         attn = attn.softmax(dim=-1)
1141 56         attn = self.attn_drop(attn)
1142 57         x = attn @ v
1143 58
1144 59         x = x.transpose(1, 2).reshape(B, N, C)
1145 60         x = self.proj(x)
1146 61         x = self.proj_drop(x)
1147 62         return x

```

H DISCUSSION ABOUT INITIALIZATION IMPLEMENTATION IN NANO GPT

We observe that, in some popular open-sourced project, *e.g.*, nanoGPT, they use an initialization implementation as code below. Let us consider a model with hidden dimension 768. Suppose

we have a linear layer projecting a 768-d feature into a new 768-d feature. For such a linear layer, the used standard variance is $\text{math.sqrt}(\frac{2}{768+768}) \approx 0.036$. For *StableNorm*, the used standard variance is $\frac{1}{2*\text{math.sqrt}(768)} \approx 0.018$. In the following code, the used standard variance is 0.02. It works. However, when we train a GPT-3 175B model with hidden dimension 12288, the standard variance 0.02 is too large. For a GPT-3 175B model with hidden dimension 12288, For *StableNorm*, the used standard variance is $\frac{1}{2*\text{math.sqrt}(12768)} \approx 0.0045$.

CODE 2: Initilization Implementation in nanoGPT.

```

1 def _init_weights(self, module):
2     if isinstance(module, nn.Linear):
3         torch.nn.init.normal_(module.weight, mean=0.0, std=0.02)
4         if module.bias is not None:
5             torch.nn.init.zeros_(module.bias)
6     elif isinstance(module, nn.Embedding):
7         torch.nn.init.normal_(module.weight, mean=0.0, std=0.02)

```

In conclusion, this implementation works for small model, but it will make training unstable or harder to train when the model is large, e.g., GPT-3 13B or GPT-3 175B.

I ABLATION STUDY

StableGPT can tolerate larger learning rate. To further validate the stability of our algorithm, we used larger learning rates (1.2e-3, 1.8e-3, 2.4e-3) to test our model. As shown in Figure 6, we found that our model can tolerate higher learning rates while maintaining good stability. Meanwhile, we can see that StableGPT-S using 1.2e-3 learning rate achieves a better performance than 6e-4 (2.819 verse 2.827).

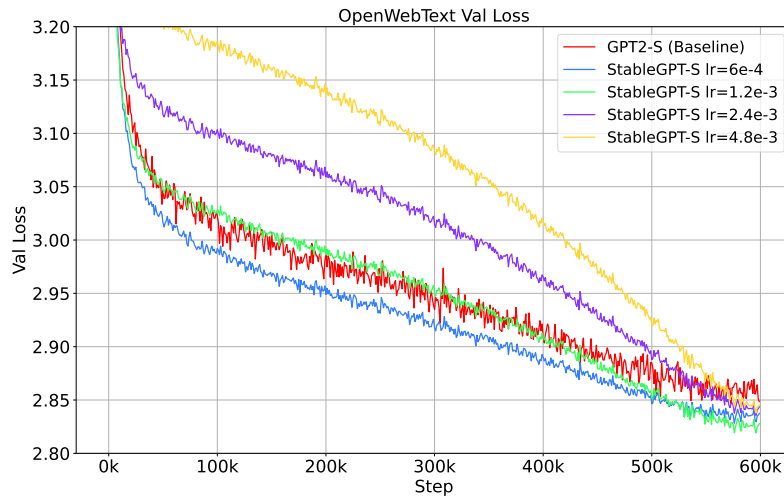


FIGURE 6: StableGPT can tolerate larger learning rate.

StableGPT is robust to the temperature coefficient in StableAtten. We conducted an evaluation of the parameter τ in the *StableAtten*, using values of $\tau = 0.809 \log_2 N$, $\tau = 1.618 \log_2 N$, and $\tau = 3.236 \log_2 N$, the used learning rate here is 6e-4 for all comparisons. We found that our algorithm is relatively robust to this parameter, with performance remaining stable across these values.

1188
 1189
 1190
 1191
 1192
 1193
 1194
 1195
 1196
 1197
 1198
 1199
 1200
 1201
 1202
 1203
 1204
 1205
 1206
 1207
 1208
 1209
 1210
 1211
 1212
 1213
 1214
 1215
 1216
 1217
 1218
 1219
 1220
 1221
 1222
 1223
 1224
 1225
 1226
 1227
 1228
 1229
 1230
 1231
 1232
 1233
 1234
 1235
 1236
 1237
 1238
 1239
 1240
 1241

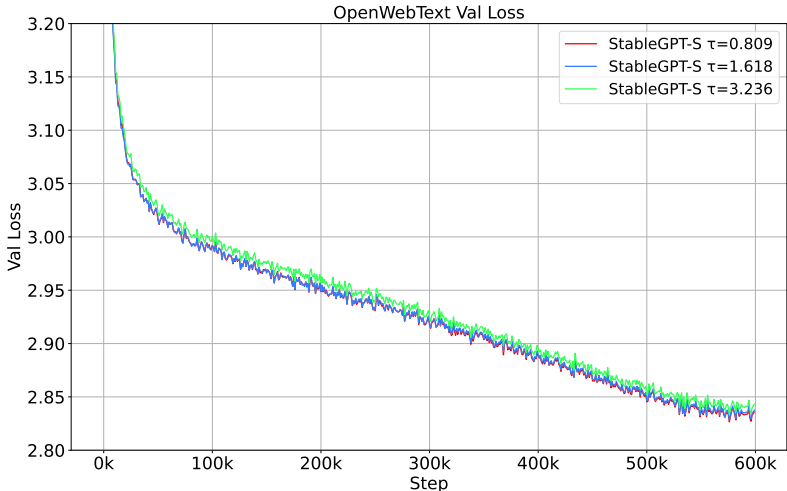


FIGURE 7: Evaluation of temperature coefficient in *StableAtten*.

About *StableViT-Huge*. We also conducted evaluations and comparisons on larger *StableViT* models, as shown in Figure 8. Compared to ViT-Huge, our algorithm demonstrates better performance, 81.8 (*StableViT-Huge*) versus 80.5 (ViT-Huge). We also noticed that Model *StableViT-Huge* is not as good as Model *StableViT-Large*, which may be mainly due to two aspects: 1). Insufficient data leading to a certain degree of overfitting, 2). Inadequate data augmentation, although we have adopted data augmentation methods similar to those in previous papers (Xie et al., 2024).

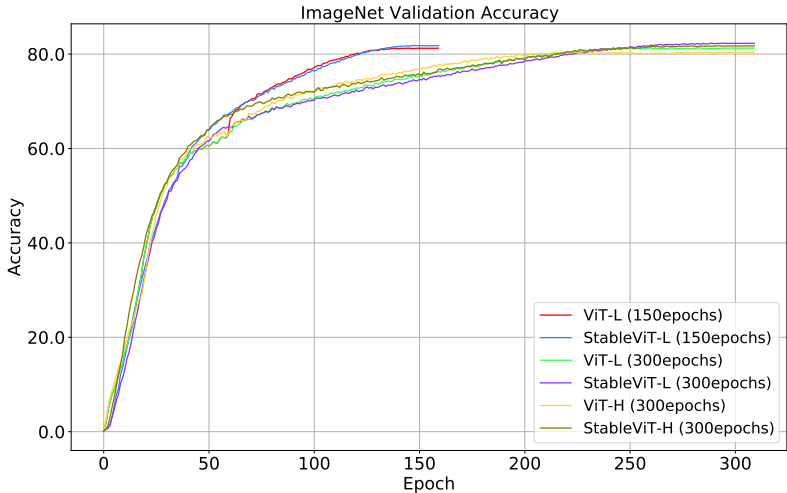


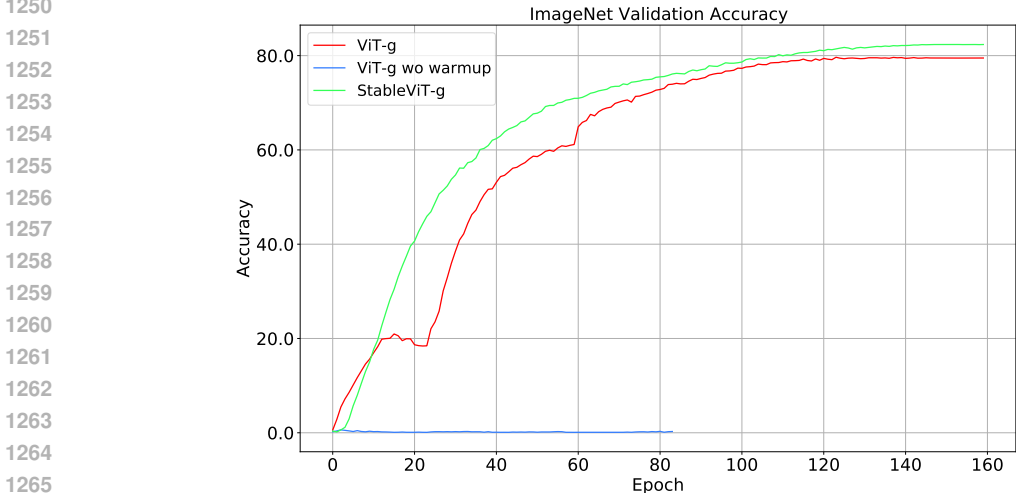
FIGURE 8: Evaluation of *StableViT-Huge*.

J EXPERIMENT OF 1B STABLEViT

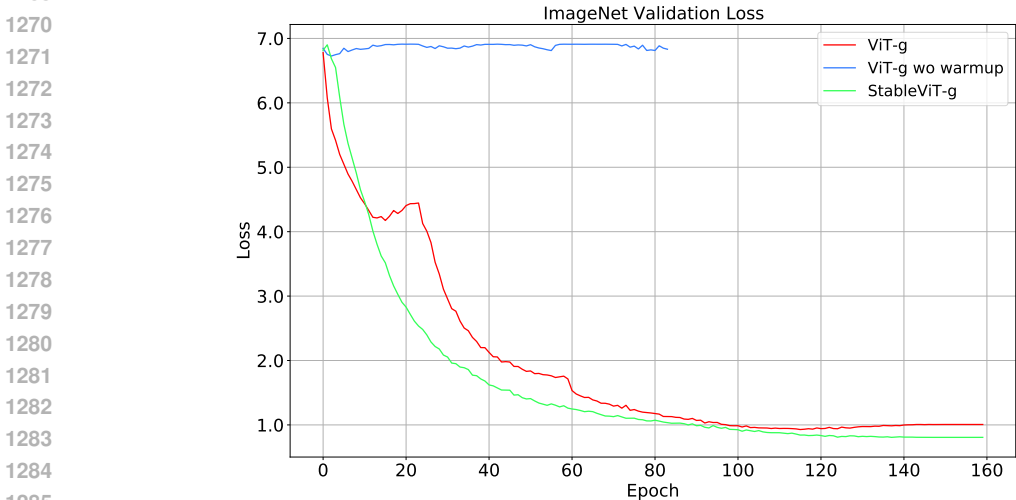
To further evaluate the effectiveness of our method at a larger scale, we assessed *StableViT* with 1B parameters, we term it as *StableViT-g* where “g” means giant. The model architecture consists of 40 layers with a hidden dimension of 1408, 16 attention heads, and an MLP dimension of 6144. The total parameter count is 1011M, around one billion parameters. We conducted a comparative study between *StableViT-g* and ViT-g, where ViT-g was evaluated under two settings: with

1242 and without learning rate warmup. Our StableViT-g does not use warmup. In StableViT, we use
 1243 an α value of 0.25. The comparison results are presented in Figure 9 and Figure 10.
 1244

1245 Figure 9 shows that ViT-g crashes after only a few training steps when running without warmup.
 1246 While the use of warmup enables ViT-g to complete training, our StableViT-g not only achieves
 1247 stable training without warmup but also demonstrates superior performance. Meanwhile, from
 1248 Figure 10, we can also observe that the loss of StableViT-g has no spike, but ViT-g even with learn-
 1249 ing rate warmup has a spike.



1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267 **FIGURE 9: Accuracy of *StableViT-g* compared with ViT-g.**



1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286 **FIGURE 10: Loss curve of *StableViT-g* compared with ViT-g.**

1287
1288
1289 **K EXPERIMENT OF 0.77B STABLEGPT**
 1290

1291 We also evaluated the effectiveness of StableGPT at a larger scale, termed as StableGPT-large.
 1292 The model architecture consists of 36 layers with a hidden dimension of 1280 and 20 attention
 1293 heads. The total parameter count is 774M. Our experimental setup strictly follows the nanoGPT
 1294 configuration, including all learning rate settings. It is important to note that training StableGPT-
 1295 large is computationally intensive, requiring two weeks to train 600K steps on 16 A800 GPUs.
 To reduce the training time, we limited our training to 100K steps instead of the full 600K steps.

The comparison results are presented in Figure 11. We can see from Figure 11, StableGPT-large obtains a better validation loss, 2.523 versus 2.536, than its counterpart, GPT2-large.

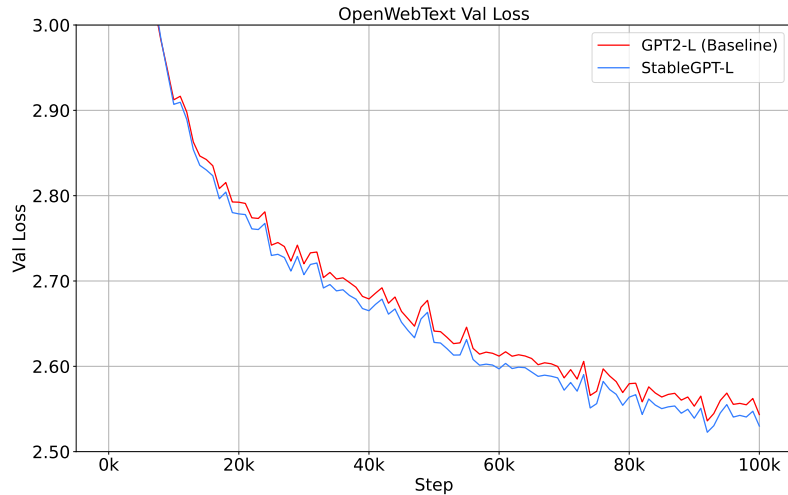


FIGURE 11: Evaluation of *StableGPT*-Large compared with *CPT2*-Large.

L EXPERIMENT OF 200 LAYERS' STABLEViT WITH 1.44B PARAMETERS

To further verify the stability of our Stable-Transformer, we conduct an experiment of super deep StableViT that has 200 layers. The model architecture consists of 200 layers with a hidden dimension of 768, 12 attention heads, and an MLP dimension of 3072. The total parameter count is 1439M, around 1.4B. We term our model as StableViT-200. Finally, StableViT-200 has 1.44B parameters. The α in StableNorm is set to be 0.25. We compare StableViT-200 with ViT-200.

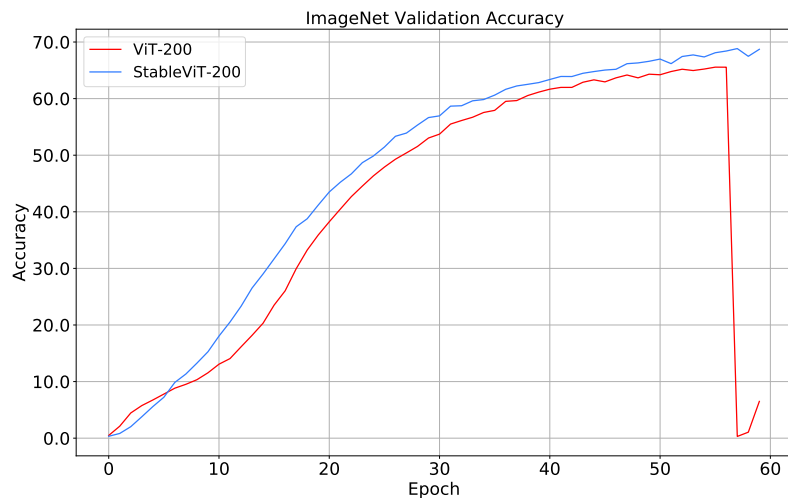


FIGURE 12: Accuracy of *StableGPT*-200 compared with ViT-200.

1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403

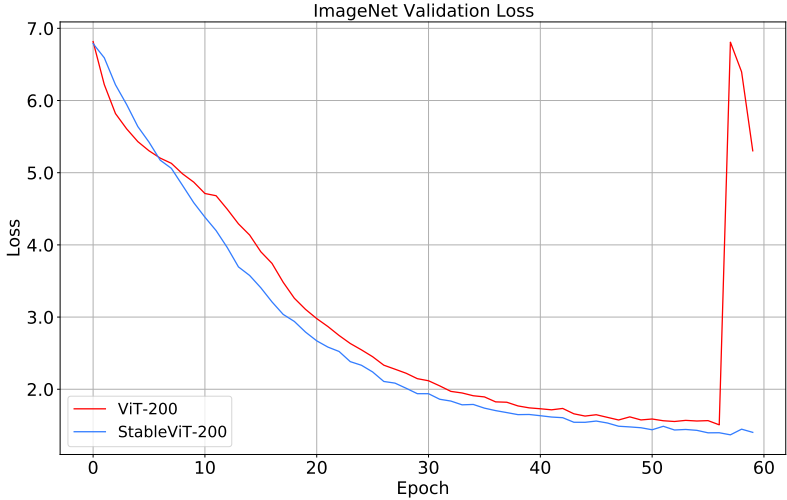


FIGURE 13: Loss curve of *StableGPT-200* compared with ViT-200.

From Figure 13, with learning rate warmup, ViT-200 has a smoothing loss curve in the early stage, but the loss spikes around at 56th epochs. But StableViT-200, without using learning rate warmup, can converge stably. It fully verify the stability of StableViT in very deep Transformer.

M DISCUSSION ABOUT LIPSCHITZ CONSTANT OF STABLENORM

Lipschitz continuity of the network is a very important condition for a stable training. Actually the principle behind StableNorm can also be explained by its Lipschitz constant. Note that the Jacobian matrix of StableNorm is defined as

$$\frac{\partial \text{StableNorm}(x)}{\partial x} = \frac{d^\alpha}{\sqrt{\|x\|_2^2 + \epsilon}} \left(\mathbf{I} - \frac{xx^\top}{\|x\|_2^2 + \epsilon} \right) \text{diag}(\gamma)$$

and the Jacobian matrix of RMSNorm is defined as

$$\frac{\partial \text{RMSNorm}(x)}{\partial x} = \frac{d^{0.5}}{\sqrt{\|x\|_2^2 + \epsilon}} \left(\mathbf{I} - \frac{xx^\top}{\|x\|_2^2 + \epsilon} \right) \text{diag}(\gamma).$$

By choosing a smaller α , e.g., $\alpha < 0.5$, the Lipschitz constant of StableNorm will less than that of RMSNorm. For example, if $d = 1024$, when we choose $\alpha = 0.475$, the Lipschitz constant of StableNorm is only around 84% of that of RMSNorm. This explains why StableNorm has a better stability than RMSNorm.

N STABLEATTEN COMPARED WITH L_2 SELF-ATTENTION

We further compared our StableAtten with L_2 self-attention (Kim et al., 2021). As shown in (Kim et al., 2021), a necessary condition to guarantee its Lipschitz continuity is $\mathbf{W}_q = \mathbf{W}_k$, thus we evaluate two versions of L_2 self-attentions: a) using tied \mathbf{W}_q and \mathbf{W}_k , i.e., $\mathbf{W}_q = \mathbf{W}_k$ and b) using two separate \mathbf{W}_q and \mathbf{W}_k , i.e., $\mathbf{W}_q \neq \mathbf{W}_k$. We conduct experiments to compare the two versions of L_2 self-attention methods with StableGPT-large, where the same training settings as the experiments in Appendix K are used, and show in Figure 14 the validation losses of our StableGPT-g and ViT-g with L_2 self-attention.

We can see from Figure 14 that, StableGPT-L with StableAtten achieves better validation loss than that of using the L_2 self-attention methods. Note that the performance degenerates notably when using the L_2 self-attention with tied \mathbf{W}_q and \mathbf{W}_k .

1404
1405
1406
1407
1408
1409
1410
1411
1412
1413
1414
1415
1416
1417
1418
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457

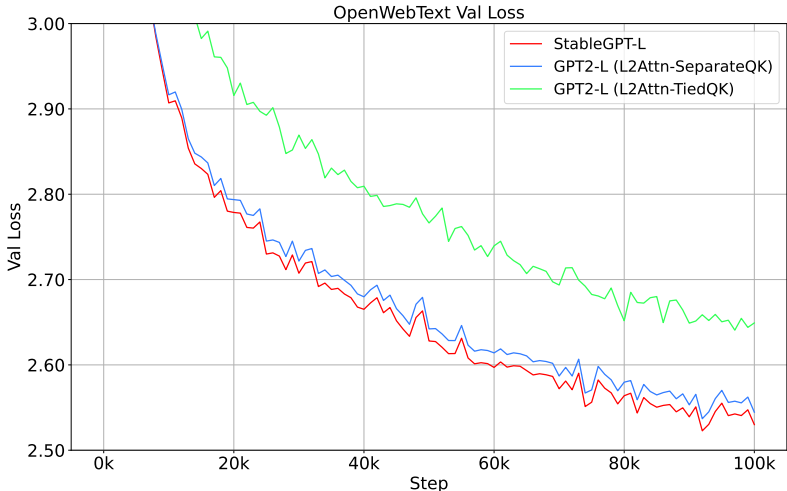


FIGURE 14: The curve of validation loss of *StableGPT-g* compared to *ViT-g* with L_2 self-attention (Kim et al., 2021) under two settings.

O ROBUSTNESS TO DISTRIBUTION SHIFT

We further conduct a set of experiments to compare the robustness between *StableViT-small* and *ViT-small* against to the distribution shift on CIFAR-100. The protocol in experiments is to train both models on the original dataset CIFAR-100 for 200 epochs, with a batch size 512, a learning rate $1e-3$, and a weight decay of $1e-4$, and then to evaluate the trained models on the original test images of CIFAR-100 and the corrupted test images of CIFAR-100, respectively. Experimental results are reported in Table 4.

TABLE 4: Evaluation (Accuracy) of robustness of *StableViT* against to distribution shift.

Models	CIFAR-100	CIFAR-100-C
<i>ViT-small</i>	67.3	51.5
<i>StableViT-small</i>	69.9	53.4

We can see that from Table 4, *StableViT-small* obtains a better accuracy than *ViT-small*, and improves the accuracy from 67.3% to 69.9%. On the corrupted CIFAR-100-C dataset, *StableViT-small* also shows a better robustness to corruption from 51.5% to 53.4%.

P RELATED WORK

Initialization. Xavier Initialization does the most groundbreaking work in model Initialization. it sets the weights to ensure the variance of activations remains constant across layers, relieving the vanishing and exploding gradient problems. Sutskever et al. (2013) investigates the importance of initialization and momentum (Nesterov, 1983; 1998) in deep learning. Kaiming Initialization (He et al., 2015), builds on Xavier Initialization by scaling the weights for ReLU activations (Nair & Hinton, 2010). Admin (Liu et al., 2020) introduces an adaptive initialization method that dynamically adjusts the initialization parameters based on the network’s depth and width. Saxe et al. (2013) introduce an orthogonal initialization, which further optimizes the initial parameter distribution to boost training outcomes. Arpit et al. (2019) also investigates the orthogonal initialization. Huang et al. (2020) propose to scale decoder by $(9L)^{-\frac{1}{4}}$ and scale encoder by $0.67L^{-\frac{1}{4}}$,

1458 this initialization method can be seen as a depth-aware initialization. *Different from the above-*
1459 *mentioned methods, our StableInit is built on Random Matrix Theory, can promise the weight*
1460 *initialized by StableInit has Lipschitz constant approximately 1.*

1461 **Normalization.** LayerNorm (Ba et al., 2016), different from BatchNorm (Ioffe & Szegedy,
1462 2015), normalizes across the features for each data point, making it effective for recurrent and
1463 transformer-based architectures. Wang et al. (2019) discuss the influence of Pre-Norm and Post-
1464 Norm on the training deep transformer. Xiong et al. (2020) further discuss the influence of
1465 pre-norm and post-norm on the training stability. RMSNorm (Zhang & Sennrich, 2019) is a
1466 variant of LayerNorm that uses root mean square statistics, offering computational efficiency.
1467 DeepNorm (Wang et al., 2022) extends normalization strategies to deep transformer networks.
1468 WeightNorm (Salimans & Kingma, 2016) reparameterizes weight vectors to decouple the magni-
1469 tude from the direction, facilitating smoother optimization. CenterNorm (Qi et al., 2023b) only
1470 conducts the centering but does not scaling the feature. ScaleNorm (Nguyen & Salazar, 2019)
1471 normalizes only by the scale of the feature vectors, simplifying the normalization process. *RM-*
1472 *SNorm and ScaleNorm can be seen as a special case of our StableNorm where $\alpha = 0.5$ and $\alpha = 0$.*
1473 *By choosing a better α , our StableNorm can obtain a better training stability.*

1474 **Attention.** Attention mechanism (Bahdanau et al., 2014) is firstly introduced to neural machine
1475 translation. Scaled dot-product attention, used in the Transformer architecture, calculates the at-
1476 tention weights using the scaled dot-product of query and key vectors, providing an efficient way
1477 to capture dependencies. L2 distance attention employs the Euclidean distance between queries
1478 and keys to compute attention scores. Attention with QK-Norm (Henry et al., 2020) normalizes
1479 the query and key vectors before computing attention, improving stability and performance. De-
1480 ghani et al. (2023) scale the model to 22B via bringing QKNorm into attention. Wortsman et al.
1481 (2024) further experimentally evaluate the value of QKNorm on small-scale models. However,
1482 these three papers do not mathematically explain why QKNorm works. Liu et al. (2022) intro-
1483 duce to use a Scaled Cosine Attention (SCA) for Transformer. Meanwhile, Qi et al. (2023a) also
1484 propose to use scaled cosine similarity attention (SCSA) to compute attention weights. Different
1485 from Liu et al. (2022), SCSA (Qi et al., 2023a) multiply a temperature coefficient instead of divid-
1486 ing a temperature coefficient. Cosine similarity attention and attention with QK-Norm share the
1487 similar idea, except that the former uses a scalar as a scale, but the latter uses a vector γ , SCSA
1488 also normalizes the values but the latter does not. *StableAtten, the logit of the attention will not be*
directly related to the hidden dimension d , and thus it is robust to the increase of the model scale.

1489 **Neural Network Stability.** To obtain a better training stability, ReZero (Bachlechner et al., 2021)
1490 introduces a simple yet effective mechanism where residual connections start as zero, allowing
1491 networks to learn identity mappings more easily and stabilize training. Admin (Liu et al., 2020)
1492 not only offers an initialization scheme but also contributes to network stability by dynamically
1493 adjusting learning rates and weight decay. DeepNorm (Wang et al., 2022) extends its benefits to
1494 network stability by adjusting normalization parameters dynamically to accommodate deeper
1495 networks. Lipsformer (Qi et al., 2023a) introduce a Lipschitz continuity constraint to ensure sta-
1496 bility in transformer networks, addressing the issue of exploding gradients. Large et al. (2024)
1497 introduces a modular norm strategy for scalable optimization. The modular norm normalizes
1498 the weights and their updates in the forward and the backward individually. They prove that
1499 the gradient of the network is Lipschitz-continuous in the modular norm with the Lipschitz con-
1500 stant that admits a simple recursive formula. The modular norm introduces a new possible di-
1501 rection for future deep neural network optimization. However, a problem is that it cannot be
1502 directly plugged into current Transformer framework. All components in Transformer needs to
1503 be re-adapted. *Our Stable-Transformer is built on our stabilized components, i.e., StableInit, Sta-*
bleNorm and StableAtten. It roots on solid theoretical justification.

1504 Some other great works also investigate the feature learning or representation learning (Yang,
1505 2019; Yang & Hu, 2021; Yang et al., 2022) and learning stability (Bernstein et al., 2020), we would
1506 like to recommend them to the readers although they are not directly related to this paper.
1507
1508
1509
1510
1511