

TAMING TRANSFORMER

WITHOUT USING LEARNING RATE WARMUP

Xianbiao Qi¹, Yelin He¹, Jiaquan Ye¹, Chun-Guang Li^{2†}, Bojia Zi³, Xili Dai⁴, Qin Zou⁵, Rong Xiao^{1†}

¹Intellifusion Inc. ²BUPT ³CUHK ⁴HKUST (GZ) ⁵WHU

ABSTRACT

Scaling Transformer to a large scale without using some technical tricks such as learning rate warmup and using an obviously lower learning rate is an extremely challenging task, and is increasingly gaining more attention. In this paper, we provide a theoretical analysis for the process of training Transformer and reveal the rationale behind the model crash phenomenon in the training process, termed *spectral energy concentration* of $\mathbf{W}_q^\top \mathbf{W}_k$, which is the reason for a malignant entropy collapse, where \mathbf{W}_q and \mathbf{W}_k are the projection matrices for the query and the key in Transformer, respectively. To remedy this problem, motivated by *Weyl's Inequality*, we present a novel optimization strategy, *i.e.*, making the weight updating in successive steps smooth—if the ratio $\frac{\sigma_1(\nabla \mathbf{W}_t)}{\sigma_1(\mathbf{W}_{t-1})}$ is larger than a threshold, we will automatically bound the learning rate to a weighted multiple of $\frac{\sigma_1(\mathbf{W}_{t-1})}{\sigma_1(\nabla \mathbf{W}_t)}$, where $\nabla \mathbf{W}_t$ is the updating quantity in step t . Such an optimization strategy can prevent spectral energy concentration to only a few directions, and thus can avoid malignant entropy collapse which will trigger the model crash. We conduct extensive experiments using ViT, Swin-Transformer and GPT, showing that our optimization strategy can effectively and stably train these Transformers without using learning rate warmup.

“Nothing in life is to be feared. It is only to be understood.”

— Marie Curie

1 INTRODUCTION

Transformer (Vaswani et al., 2017) has revolutionized various domains of artificial intelligence, including natural language processing (Radford et al., 2018; 2019; Brown et al., 2020; Chowdhery et al., 2023; Touvron et al., 2023; Dubey et al., 2024) and computer vision (Dosovitskiy et al., 2020; Liu et al., 2021) and many more applications (Radford et al., 2021; Ramesh et al., 2021; Peebles & Xie, 2023), owing to their ability to capture long-range dependencies through self-attention mechanisms. However, despite their widespread application and empirical success, training deep Transformer models remains quite challenging. Practitioners frequently encounter variant issues, such as gradient explosion (Qi et al., 2023b), rank collapse (Dong et al., 2021), entropy collapse (Zhai et al., 2023) and general training instability (Kim et al., 2021; Qi et al., 2023b), especially during the initial training stage.

To address these challenges, researchers have proposed various modifications to the original Transformer architecture, including altering the placement of Layer Normalization (Wang et al., 2019; Xiong et al., 2020) (*e.g.*, pre-LN vs. post-LN schemes), carefully conditioning the residual connections (Bachlechner et al., 2021), and QKNorm (Henry et al., 2020; Dehghani et al., 2023) for the self-attention module. Similarly, DeepNet (Wang et al., 2022) introduces a new normalization function to modify the residual connection in Transformer. ReZero (Bachlechner et al., 2021) introduces a learnable residual scalar parameter for the residual shortcut, and requires initiating it to 0 at the start stage of training. More recent approaches (Kim et al., 2021; Qi et al., 2023a) have focused on examining and enforcing Lipschitz continuity properties of Transformer

† Corresponding authors

components, which can provide insights into the network behavior and the training stability. *Although there are a few works (Bachlechner et al., 2021; Qi et al., 2023a) that can avoid using learning rate warmup to train Transformer successfully, all of them require significant modifications of the network structure.*

Learning rate warmup (Loshchilov & Hutter, 2016) seems to be a must-have technology for standard optimizers (Robbins & Monro, 1951; Duchi et al., 2011; Kingma & Ba, 2014; Loshchilov & Hutter, 2019) in some popular large Transformer models (Radford et al., 2018; 2019; Brown et al., 2020; Chowdhery et al., 2023; Touvron et al., 2023). Without the learning rate warmup stage, the Transformer training will prone to diverge.

Although it is usual to train a Transformer by modifying the network structure as mentioned above or using the learning rate warmup, two natural and interesting questions remain:

1. *What are the training dynamics of a Transformer model when its training fails or succeeds?*
2. *Can we successfully tame a Transformer without changing the network structure or without using learning rate warmup?*

This paper aims to answer these questions. To answer the first question, we examine the training processes of three types of Transformers, by visualizing the changing trajectories along with the training process of 15 (or 13) quantities about the parameters, activations, and attention maps. By doing so, we observe that the model crash is accompanied by a weird phenomenon that the entropy of the attention map is almost 0 and the spectral norm of $\mathbf{W}_q^\top \mathbf{W}_k$ increases to a very large value. By conducting mathematical analysis for the Transformer training, we identify that the Spectral Energy Concentration (SEC) of $\mathbf{W}_q^\top \mathbf{W}_k$ is the key problem leading to the model crash. To answer the second question, motivated by Weyl' Inequality, we present a novel optimization strategy, *i.e.*, making weight updating smooth, and verify empirically that our optimization strategy can prevent spectral energy concentration and thus achieving a stable convergence in training.

Paper Contributions. The contributions of the paper are highlighted as follows.

- We visualize the training dynamics of Transformers that train successfully or unsuccessfully and summarize two important observations from unsuccessful training that: a) the rank of the attention map matrix tends to very low and the entropy of attention probability matrix tends to 0; and b) $\sigma_1(\mathbf{W}_q^\top \mathbf{W}_k)$ increases rapidly to a very large value.
- We present theoretical analysis for the Transformer training, finding that the Jacobian matrix $\frac{\partial \text{vec}(\mathbf{P})}{\partial \text{vec}(\mathbf{W}_q^\top \mathbf{W}_k)} = \mathbf{X}^\top \otimes \mathbf{X}^\top$, where $\mathbf{P} = \mathbf{X}^\top \mathbf{W}_q^\top \mathbf{W}_k \mathbf{X}$. It implies that the gradient of $\mathbf{W}_q^\top \mathbf{W}_k$ is largely dominated by the rank of $\mathbf{X}^\top \otimes \mathbf{X}^\top$.
- We reveal that SEC of $\mathbf{W}_q^\top \mathbf{W}_k$ makes the attention map matrix to be sparse yet low-rank and it is the reason leading to model crash.
- Motivated by the Weyl's inequality, we introduce a novel strategy to address the problem of SEC of $\mathbf{W}_q^\top \mathbf{W}_k$ by controlling the rapid growth of singular values, and verify that our strategy leads to a stable training process.

2 PRELIMINARIES

Matrix norm. Given a matrix \mathbf{W} , its ℓ_p -norm is defined as: $\|\mathbf{W}\|_p = \sup_{\mathbf{x} \neq 0} \frac{\|\mathbf{W}\mathbf{x}\|_p}{\|\mathbf{x}\|_p}$. When $p = 2$, the induced matrix norm is the *spectral norm*, which is defined as the largest singular value of \mathbf{W} and is also expressed as the square root of the largest eigenvalue of the Gram matrix $\mathbf{W}^\top \mathbf{W}$. The spectral norm of a matrix \mathbf{W} can be calculated as: $\|\mathbf{W}\|_2 = \max_{\mathbf{x} \in S^{n-1}} \|\mathbf{A}\mathbf{x}\|_2 = \sqrt{\lambda_{\max}(\mathbf{W}^\top \mathbf{W})} = \sigma_1(\mathbf{W})$, where $\sigma_1(\mathbf{W})$ denotes the largest singular value of matrix \mathbf{W} and $\lambda_{\max}(\mathbf{W}^\top \mathbf{W})$ denotes the largest eigenvalue of $\mathbf{W}^\top \mathbf{W}$, and S^{n-1} denotes a unit sphere in \mathcal{R}^n .

Power iteration to compute matrix spectral norm. The power iteration algorithm starts with a vector \mathbf{x}_0 of unit ℓ_2 -norm. The entire iteration process is as follows: $\mathbf{x}_{k+1} = \frac{\mathbf{W}\mathbf{x}_k}{\|\mathbf{W}\mathbf{x}_k\|_2}$ for $k = 0, \dots, K-1$. At every iteration, \mathbf{x}_k is multiplied by matrix \mathbf{W} and normalized. After K iterations,

$\|x_K\|_2$ is used as the estimated spectral norm. Usually, it takes 3 to 5 iterations to converge, and thus the computation cost is cheap.

Adam Optimizer. Adam (Kingma & Ba, 2014) is currently the most widely used optimizer for training neural networks, owing to its efficiency and effectiveness. Adam can be simply defined as: $M_t = \beta_1 M_{t-1} + (1 - \beta_1) G_t$, $V_t = \beta_2 V_{t-1} + (1 - \beta_2) G_t^2$, $W_t = W_{t-1} - \alpha_t M_t \oslash \sqrt{V_t + \epsilon}$, where G_t is the gradient at step t , G_t^2 is the element-wise square of G_t , and \oslash denotes element-wise division, α_t is the learning rate at step t , and β_1, β_2 are the first-order and the second-order momentum factors, respectively.

3 TAMING TRANSFORMER REQUIRES TO REVISIT THE TRAINING DYNAMICS

To begin with, we first give some basic notions in a Transformer, which includes an attention module, an FFN module and two normalization modules that are used before the attention module and the FFN module. For the attention module, we usually use a multi-head attention which allows the model to jointly attend to information from different representations with different heads. Here, for the convenience, without losing generality, we only use a single-head attention. To be precise, we define each of them as follows:

$$\begin{aligned} \text{Attn}(X; W_q, W_k, W_v, W_o) &= W_o W_v X \text{softmax} \left(\frac{X^\top W_q^\top W_k X}{\sqrt{d_q}} \right), \\ \text{FFN}(x; W_1, W_2) &= W_2 \text{ReLU}(W_1 x), \\ \text{LN}(x) &= \gamma \odot z + \beta, \quad \text{where } z = \frac{y}{\text{std}(y)} \text{ and } y = \left(I - \frac{1}{D} \mathbf{1}\mathbf{1}^\top \right) x, \end{aligned}$$

where $X \in \mathcal{R}^{d \times n}$, $W_q \in \mathcal{R}^{d_q \times d}$, $W_k \in \mathcal{R}^{d_q \times d}$, $W_v \in \mathcal{R}^{d_v \times d}$, $W_o \in \mathcal{R}^{d \times d}$, $W_1 \in \mathcal{R}^{4d \times d}$, $W_2 \in \mathcal{R}^{d \times 4d}$, $\gamma \in \mathcal{R}^d$ and $\beta \in \mathcal{R}^d$. Note that only in a single-head definition, W_o can be put before W_v ; otherwise, it should be after a concatenation operator. For the convenience of discussion, we define:

$$A = \text{softmax} \left(\frac{P}{\sqrt{d_q}} \right) \quad \text{and} \quad P = X^\top W_q^\top W_k X,$$

where A is called the attention map and $\frac{P}{\sqrt{d_q}}$ is usually called the logit.

3.1 VISUALIZATION: WHAT HAPPENS WHEN A TRANSFORMER TRAINING FAILS OR SUCCEEDS

Visualizations are commonly used as an effective means to help us examine whether the neural network's training succeeds or fails. In particular, what happens when the training of a Transformer tends to crash? And what happens when the training succeeds?

One of the most important aspects of understanding the training of neural networks is the observation of changes in parameters and activations. Since the parameters or activations and their gradients are matrices or vectors, the norm is the best way to observe these quantities. In this paper, for a Transformer training, we summarize the following **15 terms** to watch:

$$\begin{aligned} &\sigma_1(W_q), \quad \sigma_1(W_k), \quad \sigma_1(W_v), \quad \sigma_1(W_o), \quad \sigma_1(W_1), \quad \sigma_1(W_2), \quad \sigma_1(W_q^\top W_k), \\ &\sigma_1(W_o W_v), \quad \|\gamma_1\|_2, \quad \|\beta_1\|_2, \quad \|\gamma_2\|_2, \quad \|\beta_2\|_2, \quad \sigma_1(W_2 W_1), \quad \|x\|_2, \quad \left\| \frac{\partial L}{\partial x} \right\|_2, \end{aligned} \quad (1)$$

where β_1 and γ_1 denote the parameters of the first LayerNorm, and β_2 and γ_2 denote the parameters of the second LayerNorm. When RMSNorm (Zhang & Sennrich, 2019) is used, there are only 13 terms that will be analyzed since it does not have β . For the weight matrix, we use the spectral norm. For a vector, we use its ℓ_2 norm.

To ensure that the phenomena we observed can generalize well, we visualized them on both ViT (Dosovitskiy et al., 2020) and GPT (Radford et al., 2018). ViT is a pure encoder architectures whereas GPT is a pure decoder architecture.

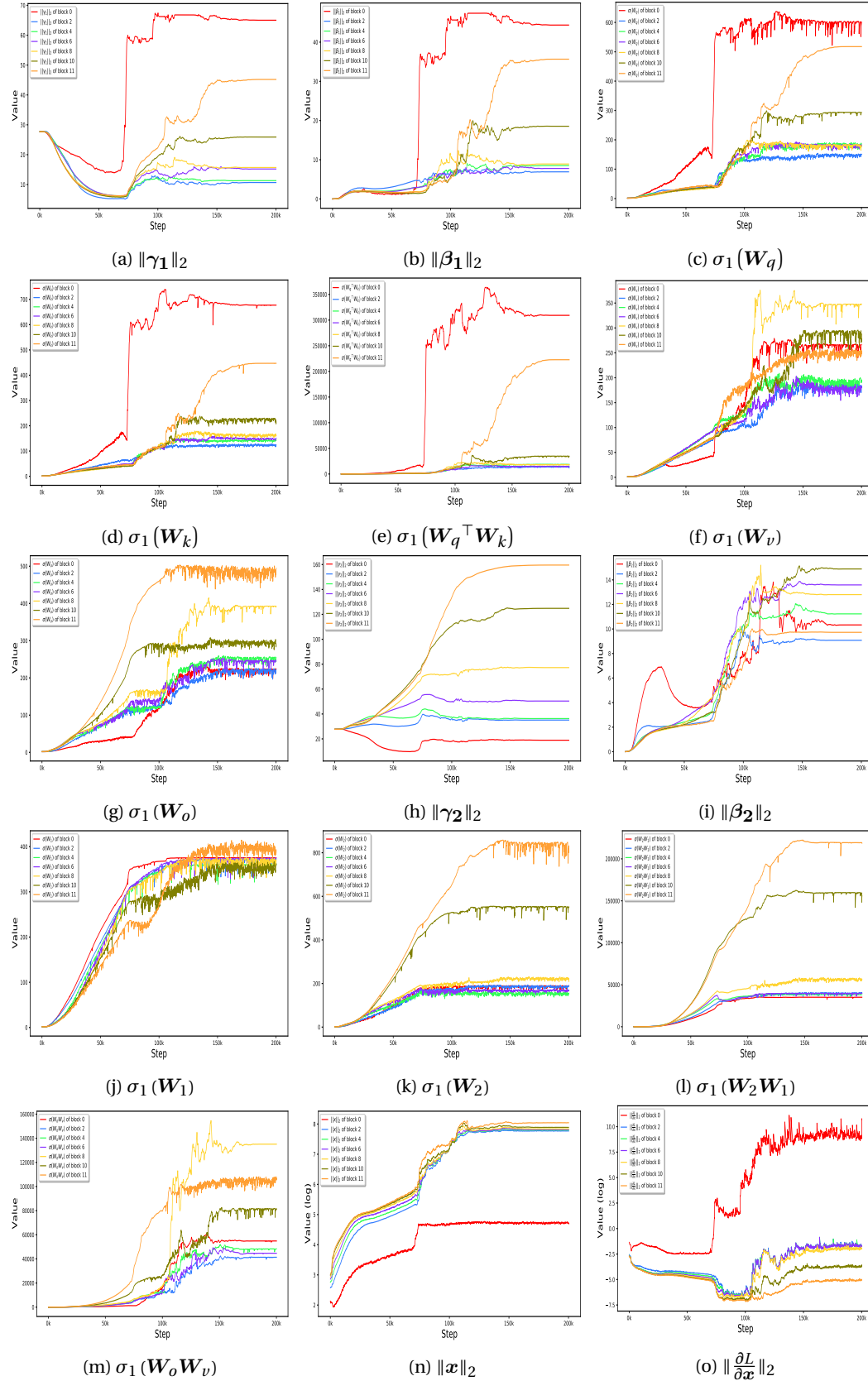


FIGURE 1: Training dynamics of a failure ViT. This figure shows how 15 items of quantities as defined in Equation (1) change during the training period. Please pay more attention to subfigures (a)-(e).

(a) Block 0 (successful). (b) Block 6 (successful). (c) Block 11 (successful).
 (d) Block 0 (crashed). (e) Block 6 in (crashed). (f) Block 11 (crashed).

FIGURE 2: Visualization of the dynamics process of attention map in different training steps for a successful and a crashed ViT-base model, respectively. *Please click the images to play the flash. Best viewed with Acrobat Reader.*

In this section, due to the space limitation, we will only visualize ViT-base, and put more visualization results into the Appendix H. For the ViT implementation, we use Timm Wightman (2019), in which “timm” library provides rich model architectures of many pre-trained image models in PyTorch. For the GPT implementation, we use nanoGPT, which uses LayerNorm without a bias term, and thus only watch 13 terms, rather than 15 terms in ViT.

To achieve a successful ViT training, we use a long learning rate warmup. For instance, we use 60 epochs of warmup and the whole training process takes 150 epochs. To

obtain the dynamics of a failure training of ViT, we do not use warmup.

Figure 1 visualizes a failure training process of a ViT-base model. The model includes 12 blocks with index from 0 to 11. In Figure 1, we visualize the weight matrices in blocks $\{0, 2, 4, 6, 8, 10, 11\}$, where the index of the last block is 11. For the features x and the gradients $\frac{\partial L}{\partial x}$, we hook the input features that enter into the corresponding blocks. Figure 8 in the Appendix I visualizes a successful training process of a ViT-L model. Meanwhile, in Figure 2, we visualize the dynamic processes of attention maps during the training period about a successful case and a failure case of the ViT-base model, respectively. We visualize the attention map of the same image at different steps.

We observe the following phenomena from Figures 1 and 2.

- As shown in Figure 1, at the beginning, the maximum singular value $\sigma_1(\mathbf{W}_q^\top \mathbf{W}_k)$ gradually increases, and at a certain point, the maximum singular value suddenly and rapidly increases to a very large value (e.g., around 200,000), at where the loss divergence emerges. However, for a successful training process, $\sigma_1(\mathbf{W}_q^\top \mathbf{W}_k)$ gradually increases to a medium value as shown in Figure 8 and then vibrates around that value.
- As shown in Figure 2, in a failure training process of Transformer, the attention maps gradually become sparse and low-rank. Finally, the entropy of the attention map is 0; whereas in a successful training process, the attention map is not too sparse and of a medium rank.
- The normalization layers exhibit a huge difference: γ_1 and β_1 in a successful ViT training process are very smooth, but $\|\gamma_1\|_2$ and $\|\beta_1\|_2$ suddenly increase dramatically in an unsuccessful case.
- The ranges of the activation and the gradients are very large in a crashed model, and the gradients in different blocks vary much larger than that in a successful model.

Remark. We summarize that the successful training and the unsuccessful training of a Transformer exhibit significant differences among their $\sigma_1(\mathbf{W}_q^\top \mathbf{W}_k)$, their normalization parameters γ and β , and their activations x and gradients $\frac{\partial L}{\partial x}$.

3.2 THEORETICAL ANALYSIS: MATRIX CALCULUS OF TRANSFORMER

To understand the training dynamics of Transformer, we should investigate the process of back-propagation (Rumelhart et al., 1986; LeCun et al., 2002; 1989; 1998). In the attention mechanism, however, the input and the output are both matrices, we cannot directly use vector calculus. Instead, we need to use *Vectorization* (Graham, 2018; Petersen et al., 2008) and *Kronecker Product* (Graham, 2018; Petersen et al., 2008). In matrix calculus, the vectorization of a matrix is a linear transformation that converts a matrix into a vector. Specifically, the vectorization of a matrix $M \in \mathcal{R}^{m \times n}$, denoted as $\text{vec}(M)$, is a column vector, obtained by an ordered stacking of the columns of the matrix M , i.e., $\text{vec}(M) \in \mathcal{R}^{mn}$. For example, for a 2×3 matrix $M = \begin{bmatrix} a & b & c \\ d & e & f \end{bmatrix}$, the vectorization of M is $\text{vec}(M) = [a \ d \ b \ e \ c \ f]^\top$. For the attention module of Transformer, we have the following proposition about the Jacobian matrix of the output P with respect to the input X and the parameters.

Proposition 1 (Matrix Calculus for Self-Attention)

Let $P = X^\top W_q^\top W_k X$, where $X \in \mathcal{R}^{d \times n}$, $W_q \in \mathcal{R}^{d_q \times d}$, $W_k \in \mathcal{R}^{d_q \times d}$, according to vectorization and matrix calculus, we have the following derivations:

$$\frac{\partial \text{vec}(P)}{\partial \text{vec}(W_q^\top W_k)} = X^\top \otimes X^\top, \quad \frac{\partial \text{vec}(P)}{\partial \text{vec}(X)} = (X^\top W_k^\top W_q \otimes I_n)K + (I_n \otimes X^\top W_q^\top W_k), \quad (2)$$

$$\frac{\partial \text{vec}(P)}{\partial \text{vec}(W_q^\top)} = (W_k X)^\top \otimes X^\top, \quad \frac{\partial \text{vec}(P)}{\partial \text{vec}(W_k)} = X^\top \otimes (W_q X)^\top, \quad (3)$$

where \otimes denotes the Kronecker product, $I_n \in \mathcal{R}^{n \times n}$ denotes an identity matrix with shape $n \times n$, K is the commutation matrix, which depends on the dimensions of X . Since $X \in \mathcal{R}^{d \times n}$, then we know $K \in \mathcal{R}^{nd \times nd}$. The commutation matrix K has the property that $\text{vec}(X^\top) = K \text{vec}(X)$ for any matrix X .

In Appendices A and B, we supply some elementary information for the vectorization and the Kronecker product and the derivations of the Jacobian matrix for a single-head attention.

We have the following observations from Proposition 1.

- We have $\frac{\partial \text{vec}(P)}{\partial \text{vec}(W_q^\top W_k)} = X^\top \otimes X^\top$ in Equation 2, and we know about the Kronecker product that $\text{rank}(X^\top \otimes X^\top) = \text{rank}(X^\top)^2$, which implies that if the rank of X is very low, then the rank of $\frac{\partial \text{vec}(P)}{\partial \text{vec}(W_q^\top W_k)}$ will also be very low. Note that X being low rank means that the features across different timestep are highly correlated or coherent. If all x_i in X collapses to a single point, then $X^\top \otimes X^\top$ will only have a large singular value, and the rest are 0.
- The Jacobian matrix $\frac{\partial \text{vec}(P)}{\partial \text{vec}(X)}$ in Equation 2, is in direct proportion to X and $W_q^\top W_k$. If the spectral norm $\sigma_1(W_q^\top W_k)$ is very large, it implies that the gradient $\frac{\partial L}{\partial X}$ will more likely to be magnified a lot.
- Equation 3 suggests that changes in the query weights W_q are related to both the input X and the key representation $W_k X$. Equation 3 suggests that changes in the key weights W_k are related to both the input X and the query representation $W_q X$.
- All these relationships are interconnected, with changes in one variable potentially affecting the others. For instance, if W_k increases fast, then according to Equation 3, $\frac{\partial \text{vec}(P)}{\partial \text{vec}(W_q^\top)}$ will more likely to be very large. In this way, W_q will likely increase very fast.

3.3 RATIONALE IN MODEL CRASH: SPECTRAL ENERGY CONCENTRATION

Before we reveal the rationale in model crash, let us first discuss a bit on the entropy collapse.

Two Entropy Collapse Modes. In experiments, we observe two types of attention entropy collapse modes. Note that when attention collapse happens, the attention map tends to

[https://en.wikipedia.org/wiki/Vectorization_\(mathematics\)](https://en.wikipedia.org/wiki/Vectorization_(mathematics))
https://en.wikipedia.org/wiki/Kronecker_product

a sparse matrix (*i.e.*, there are a few dominate nonzero coefficients in attention map), and thus the entropy of the attention map is vanishing. To be more specific, when the attention map is sparse but not low-rank, we call it a *benign collapse*; whereas if the attention map is sparse and simultaneously low-rank, we call it a *malignant collapse*. When benign collapse occurs, the attention map is shown in the right panel of Figure 3 that, there is almost an identity matrix. In this way, the diagonal elements are almost 1, and the non-diagonal elements have values around 0. Unfortunately, when malignant collapse happens, the attention map becomes a sparse and simultaneously low-rank matrix as shown in the middle panel of Figure 3. Furthermore, we observe that the distribution of the spectral energy of $\mathbf{W}_q^\top \mathbf{W}_k$ for the benign collapse is relatively uniform; whereas the spectral energy of the attention matrix for the malignant collapse tends to concentrate on a few dominate singular values.

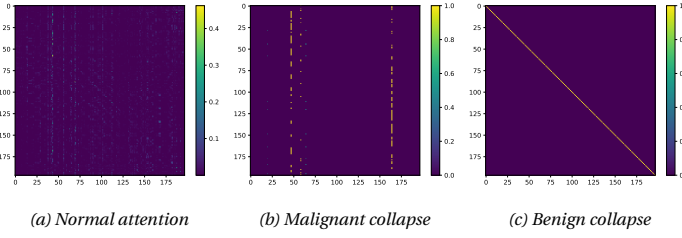


FIGURE 3: Three modes of attention maps. The left panel shows a normal attention map. The middle panel shows a classical attention map when model crash occurs, for which the entropy is almost 0. The right panel shows an attention map from a normal model training but its entropy is almost 0.

By analyzing the matrix $\mathbf{W}_q^\top \mathbf{W}_k$ in the benign collapse when it happens in the experiments, we find that it has the following property: $\mathbf{W}_q^\top \mathbf{W}_k$ is usually a non-symmetric positive quasi-definite square matrix (see Appendix O for details), and the self-attention layer degenerates into a linear projection layer because

$\mathbf{Y} = \mathbf{W}_v \mathbf{X} \mathbf{A} \approx \mathbf{W}_v \mathbf{X} \mathbf{I} = \mathbf{W}_v \mathbf{X}$. We give an intuitive analysis in Appendix C.

When the malignant collapse happens, the model will usually crash. We identify that the rationale behind the model crash is a phenomenon called *spectral energy concentration (SEC)*. Before we present our theorem about SEC, let us first introduce an index to quantify it. Recall that $\mathbf{W}_q \in \mathcal{R}^{d_q \times d}$ and $\mathbf{W}_k \in \mathcal{R}^{d_q \times d}$, where $d_q < d$. We have that $\mathbf{W}_q^\top \mathbf{W}_k \in \mathcal{R}^{d \times d}$, but its rank is less than or equal to d_q . Precisely, we define a *SEC index* as follows:

$$\text{SEC}(d_q, s) = \frac{\sum_{i=1}^s \sigma_i^2(\mathbf{W}_q^\top \mathbf{W}_k)}{\sum_{i=1}^{d_q} \sigma_i^2(\mathbf{W}_q^\top \mathbf{W}_k)}, \quad (4)$$

where d_q is the head dimension and $s \leq d_q$. For instance, if we have $d_q = 64$ and $s = 4$, and if at this time, $\text{SEC}(64, 4) > 99\%$, we could say the spectral energy of $\mathbf{W}_q^\top \mathbf{W}_k$ highly concentrates on only four dominant singular values.

To be precise, we have the following theorem for the reason to cause a malignant collapse.

Theorem 1 (Malignant Entropy Collapse)

Let $\mathbf{P} = \mathbf{X}^\top \mathbf{W} \mathbf{X}$ and $\mathbf{A} = \text{softmax}\left(\frac{\mathbf{P}}{\sqrt{d_q}}\right)$, $\mathbf{W} = \mathbf{W}_q^\top \mathbf{W}_k \in \mathbb{R}^{d \times d}$. Suppose that the following two conditions are simultaneously satisfied:

- \mathbf{X} is a low-rank matrix;
- \mathbf{W} is a low-rank matrix with only a few dominant singular values (e.g., the singular values are greater than $C_0 \cdot \sqrt{d_q}$ where $C_0 \gg 1$ is a constant).

Then, the attention map \mathbf{A} will be sparse and simultaneously low-rank in high probability.

When the malignant entropy collapse happens, the training process will crash. We provide a proof for Theorem 1 in Appendix D.

For a sparse and simultaneously low-rank matrix, we refer the reader to a recent textbook (Wright & Ma, 2022).

According to Theorem 1, we know that the model crash is caused by the spectral energy concentration. In Figure 4, we evaluate and compare the curves of $\text{SEC}(d_q, s)$ of three different blocks under a successfully trained model and a crashed model. In a successfully trained model, the spectral energy distributes in all directions. However, the spectral energy of a crashed model only concentrates on a few directions. As shown in Figure 4, in a crashed model, the SEC collapses into less than 10 directions. Figure 5 reveals how the attention collapse propagates through

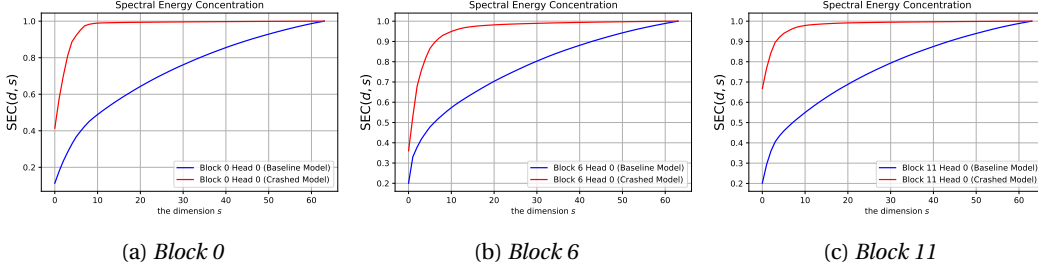


FIGURE 4: Comparison of spectral energy concentration index SEC between a successfully trained model and a crashed model. The attention maps of three different blocks are shown. The spectral energy distributes in all directions in a successful training case; whereas the spectral energy only concentrates on a few directions in a crashed model.

$$X^l \downarrow \longrightarrow X^\top \otimes X^\top \downarrow \longrightarrow \frac{\partial \text{vec}(\mathbf{P})}{\partial \text{vec}(\mathbf{W}_q^\top \mathbf{W}_k)} \downarrow \longrightarrow \mathbf{W}_q^\top \mathbf{W}_k \downarrow \longrightarrow \mathbf{A} \Downarrow \longrightarrow X^{l+1} \downarrow$$

FIGURE 5: Attribution flow chart of attention collapse.

each term, illustrating the entire attribution chain from the input X^l to the output X^{l+1} in the attention module. Note that \downarrow indicates being low-rank and \Downarrow means being sparse and simultaneously low-rank. If X is low-rank, then $X^\top \otimes X^\top$ is also low-rank because $\text{rank}(X \otimes X) = \text{rank}(X) \cdot \text{rank}(X)$. According to the gradient computation, we have $\frac{\partial \text{vec}(\mathbf{P})}{\partial \text{vec}(\mathbf{W}_q^\top \mathbf{W}_k)} = X^\top \otimes X^\top$, thus we know $\frac{\partial \text{vec}(\mathbf{P})}{\partial \text{vec}(\mathbf{W}_q^\top \mathbf{W}_k)}$ is also low-rank. Meanwhile, it should be noted that the spectral energy of $\frac{\partial \text{vec}(\mathbf{P})}{\partial \text{vec}(\mathbf{W}_q^\top \mathbf{W}_k)}$ is over-concentrated, which means that the gradient update will dramatically change $\mathbf{W}_q^\top \mathbf{W}_k$, thus $\mathbf{W}_q^\top \mathbf{W}_k$ will have large probability to be low-rank. In the paper, we have proved that being low-rank and the leading singular values of $\mathbf{W}_q^\top \mathbf{W}_k$ are very large will lead to the attention map \mathbf{A} over-concentrated (see Appendix C for the proof), and become to a sparse and simultaneously low-rank matrix. Finally, an over-concentrated \mathbf{A} will lead to X^{l+1} to be low-rank. We provide a more detailed analysis of Figure 5 in Appendix Q.

3.4 OUR SOLUTION: TAMING TRANSFORMER VIA WEYL'S INEQUALITY

The analysis above reveals that spectral energy concentration is the key reason leading to unstable training. One manifestation of spectral energy concentration is the rapid growth of the singular values of the weight matrices. Thus, our basic idea to prevent malignant collapse is to suppress the fast growth of the singular values. Fortunately, Weyl's inequality provides us a simple but effective tool.

Theorem 2 (Weyl's Inequality on Singular Values)

Let $\mathbf{W}_1, \mathbf{W}_2 \in \mathbb{R}^{m \times n}$ where $m \geq n$, $\sigma_1(\mathbf{W}_1) \geq \sigma_2(\mathbf{W}_1) \geq \dots \geq \sigma_n(\mathbf{W}_1)$ be the ordered singular values of \mathbf{W}_1 , and $\sigma_1(\mathbf{W}_2) \geq \sigma_2(\mathbf{W}_2) \geq \dots \geq \sigma_n(\mathbf{W}_2)$ be the ordered singular values of \mathbf{W}_2 . Then we have:

$$\sigma_{i+j-1}(\mathbf{W}_1 + \mathbf{W}_2) \leq \sigma_i(\mathbf{W}_1) + \sigma_j(\mathbf{W}_2).$$

The proof for Theorem 2 is provided in Appendix E. From Theorem 2, it is easy to see that $\sigma_1(\mathbf{W}_1 + \mathbf{W}_2) \leq \sigma_1(\mathbf{W}_1) + \sigma_1(\mathbf{W}_2)$. Let \mathbf{W}_t be the weight matrix at time step t , $\nabla \mathbf{W}_t$ be the quantity computed from the gradients and their derivations (where $\nabla \mathbf{W}_t$ can be obtained by

Algorithm 1 AdamW²: Taming Transformer via Weyl' Inequality without learning rate warmup.

Input: learning rate scheduler α_t , weight decay λ , and first-order and second-order momentum β_1, β_2
Output: updated weight \mathbf{w}_T

```

1: for  $t = 1, 2, \dots, T$  do
2:    $\mathbf{G}_t = \nabla \mathbf{W}_{t-1} \mathcal{L}$ 
3:    $\mathbf{M}_t = \beta_1 \mathbf{M}_{t-1} + (1 - \beta_1) \mathbf{G}_t$ ,    $\mathbf{V}_t = \beta_2 \mathbf{V}_{t-1} + (1 - \beta_2) \mathbf{G}_t^2$ 
4:    $\hat{\mathbf{M}}_t = \mathbf{M}_t / (1 - \beta_1^t)$ ,    $\hat{\mathbf{V}}_t = \mathbf{V}_t / (1 - \beta_2^t)$ 
5:    $\nabla \mathbf{W}_t = \hat{\mathbf{M}}_t \oslash \sqrt{\hat{\mathbf{V}}_t + \epsilon}$  ▷  $\nabla \mathbf{W}_t$  is the final update quantity.
6:    $\hat{\sigma}_1 = \text{PowerIter}(\nabla \mathbf{W}_t)$ ,    $\hat{\sigma}_1 = \text{PowerIter}(\mathbf{W}_{t-1})$  ▷ Power iteration to compute spectral norm.
7:   if  $\frac{\alpha_t \hat{\sigma}_1}{\hat{\sigma}_1} > \tau$  then ▷ If Rule 1 does not meet, adjust the learning rate  $\alpha_t$ .
8:      $\alpha_t = \frac{\tau \hat{\sigma}_1}{\hat{\sigma}_1}$ 
9:   end if
10:  if Weight Decay is Yes then
11:     $\mathbf{W}_t = \mathbf{W}_{t-1} - \alpha_t \nabla \mathbf{W}_t - \alpha_t \lambda \mathbf{W}_{t-1}$ 
12:  else
13:     $\mathbf{W}_t = \mathbf{W}_{t-1} - \alpha_t \nabla \mathbf{W}_t$ 
14:  end if
15: end for

```

SGD (Robbins & Monro, 1951), Adagrad (Duchi et al., 2011), or Adam (Kingma & Ba, 2014)), α_t is the learning rate at time step t . Usually, our update equation is $\mathbf{W}_t = \mathbf{W}_{t-1} - \alpha_t \nabla \mathbf{W}_t$. According to Weyl's Inequality, we have,

$$\sigma_1(\mathbf{W}_t) = \sigma_1(\mathbf{W}_{t-1} - \alpha_t \nabla \mathbf{W}_t) \leq \sigma_1(\mathbf{W}_{t-1}) + \alpha_t \sigma_1(\nabla \mathbf{W}_t). \quad (5)$$

An important observation from Equation 5 is that if $\sigma_1(\nabla \mathbf{W})$ is very large, then \mathbf{W}_t will be significantly different from \mathbf{W}_{t-1} . It means that the successive updates of \mathbf{W}_t at time step t from time step $t-1$ would “jump” too much. A “smoother” updating should satisfy the following rule.

Rule 1 (Steady Weights Updating Rule)

Given weights matrix \mathbf{W}_{t-1} and the updating quantity $\nabla \mathbf{W}_t$ at step t , with the learning rate α_t , a steady weights should satisfy the following inequality: $\|\mathbf{W}_{t-1} - \alpha_t \nabla \mathbf{W}_t\|_2 \leq (1 + \tau) \|\mathbf{W}_{t-1}\|_2$, where $\tau > 0$ is a small factor.

To meet Rule 1, what we need is that $\sigma_1(\mathbf{W}_{t-1}) + \alpha_t \sigma_1(\nabla \mathbf{W}_t) \leq (1 + \tau) \|\mathbf{W}_{t-1}\|_2 = (1 + \tau) \sigma_1(\mathbf{W}_{t-1})$. It is easy to see that:

$$\alpha_t \leq \tau \frac{\sigma_1(\mathbf{W}_{t-1})}{\sigma_1(\nabla \mathbf{W}_t)}. \quad (6)$$

This inequality tell us that the learning rate α_t should be bounded by a ratio of singular values $\sigma_1(\mathbf{W}_{t-1})$ and $\sigma_1(\nabla \mathbf{W}_t)$. Generally, τ is a small value, e.g., 0.004 or 0.005. The intuition behind is that if the spectral norm of $\nabla \mathbf{W}_t$ is significantly larger than that of \mathbf{W}_t , then the model is potentially undergoing rapid changes. In such cases, a large learning rate could lead to training instability. Therefore, our strategy is that if $\alpha_t > \tau \frac{\sigma_1(\mathbf{W}_{t-1})}{\sigma_1(\nabla \mathbf{W}_t)}$ then we truncate α_t to $\tau \frac{\sigma_1(\mathbf{W}_{t-1})}{\sigma_1(\nabla \mathbf{W}_t)}$. Since our base optimizer is AdamW (Loshchilov & Hutter, 2019) and our algorithm is motivated by Weyl's Inequality, we term our algorithm as AdamW².

For clarity, we summarize our AdamW² in Algorithm 1, where lines 6-9 highlight our improvements over the base optimizer, the other codes are same as AdamW. According to line 6 in Algorithm 1, $\sigma_1(\nabla \mathbf{W}_t)$ and $\sigma_1(\mathbf{W}_{t-1})$ are computed via a fast power iteration method. In practice, we set the maximum iterations in power iteration to 3. Actually, we find that two iterations are enough to estimate the spectral norm of matrices. According to Equation 6, if $\alpha_t > \tau \frac{\sigma_1(\mathbf{W}_{t-1})}{\sigma_1(\nabla \mathbf{W}_t)}$, then the learning rate α_t will be truncated to $\tau \frac{\sigma_1(\mathbf{W}_{t-1})}{\sigma_1(\nabla \mathbf{W}_t)}$, or the algorithm will adjust α_t and use the default learning rate set by the learning rate schedule.

4 EXPERIMENTS

Compared to some previous works (Bachlechner et al., 2021; Wang et al., 2019; Xiong et al., 2020; Wang et al., 2022; Qi et al., 2023b) that focuse on improving the training stability of Transformer,

TABLE 1: Quantitative comparison of AdamW and AdamW² with and without learning rate warmup. AdamW² demonstrates a very competitive performance compared to AdamW.

Method Configurations Parameters	ViT (Acc. \uparrow)		GPT (Loss \downarrow)		Swin-Transformer (Acc. \uparrow)	
	ViT-B 86M	ViT-L 307M	GPT-S 125M		Swin-S 50M	Swin-B 88M
AdamW (with warmup)	80.22	81.65	2.848		83.02	83.48
AdamW ² (no warmup)	80.58	81.82	2.840		83.14	83.44

our AdamW² does not need to adjust the network structure and we do not use learning rate warmup. For ViT, Swin-Transformer and GPT, we will use a warmup of 60 epochs, 20 epochs, and 2000 steps, respectively. In AdamW², we directly use a cosine learning rate schedule and decay the learning rate from maximum to minimum.

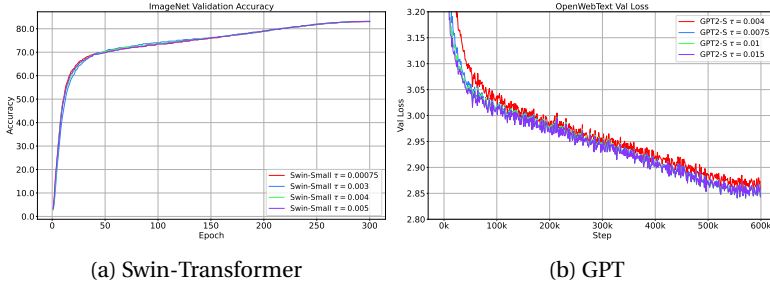


FIGURE 6: Ablation study of τ in AdamW² using Swin-S and GPT-S.

decoder. Note that we do not conduct any adjustments to the networks and directly use the original implementation. Our experiments include image classification on ImageNet (Deng et al., 2009) and large language model on OpenWebText (Gokaslan & Cohen) dataset. We list some training configurations in Appendix N. The quantitative results are shown in Table 1. Our baseline model is the corresponding Transformer using a learning rate warmup; whereas baseline models without using learning rate warmup will crash. AdamW² demonstrates a very competitive performance compared to the baseline method. These experimental results verify that our understanding of the training dynamics of the Transformer is rational.

We also conduct an ablation study of the choice of τ in GPT and Swin-Transformer. The results are shown in Figure 6. We can see that the performance of AdamW² varies slightly for different values of τ , but overall, our approach is robust for different choices of τ . The curves basically overlap in the later epochs because our steady updating rule is never broken in the later epochs.

5 CONCLUSION

In this paper, we revisited the training dynamics of the Transformers by visualizing the spectral norm of weight matrices, the activations and the attention map, presented a theoretical analysis for the Transformer training and identified two modes of attention entropy collapse, *i.e.*, the benign collapse and the malignant collapse, in which the malignant collapse accompanies model crash. Moreover, we revealed that the *spectral energy concentration* of $\mathbf{W}_q^\top \mathbf{W}_k$ is the reason behind the model crash, which causes the attention map to be sparse and simultaneously low-rank. Furthermore, we proposed a steady updating rule to resolve the problem of spectral energy concentration of $\mathbf{W}_q^\top \mathbf{W}_k$ by controlling the rapid growth of singular values, which can prevent the fast spectral energy concentration to a few directions and thus avoid the malignant entropy collapse. We conducted extensive experiments to verify the proposed strategy with ViT, Swin Transformer, and GPT, and demonstrated that the proposed strategy could effectively and stably train a model without using any learning rate warmup.

We conduct experiments on three popular Transformers, *i.e.*, ViT (Dosovitskiy et al., 2020), GPT-2 (Radford et al., 2019) and Swin-Transformer (Liu et al., 2021), where ViT and Swin-Transformer are pure encoder architectures and GPT is a pure causal

ETHICS STATEMENT

In this paper, we aim to provide a novel approach to train transformers without learning rate warmup. Our work does not involve any human subjects, and we have carefully ensured that it poses no potential risks or harms. Additionally, there are no conflicts of interest, sponsorship concerns, or issues related to discrimination, bias, or fairness associated with this study. We have taken steps to address privacy and security concerns, and all data used comply with legal and ethical standards. Our work fully adheres to research integrity principles, and no ethical concerns have arisen during the course of this study.

REPRODUCIBILITY STATEMENT

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

ACKNOWLEDGMENTS

Chun-Guang Li was partially supported by the National Natural Science Foundation of China under Grant 61876022. Qin Zou was partially funded by the National Natural Science Foundation of China under Grant 62171324.

REFERENCES

- Thomas Bachlechner, Bodhisattwa Prasad Majumder, Henry Mao, Gary Cottrell, and Julian McAuley. Rezero is all you need: Fast convergence at large depth. In *Uncertainty in Artificial Intelligence*, pp. 1352–1361. PMLR, 2021.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm: Scaling language modeling with pathways. *Journal of Machine Learning Research*, 24(240): 1–113, 2023.
- Aaron Defazio, Harsh Mehta, Konstantin Mishchenko, Ahmed Khaled, Ashok Cutkosky, et al. The road less scheduled. *arXiv preprint arXiv:2405.15682*, 2024.
- Mostafa Dehghani, Josip Djolonga, Basil Mustafa, Piotr Padlewski, Jonathan Heek, Justin Gilmer, 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.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pp. 248–255. Ieee, 2009.
- Yihe Dong, Jean-Baptiste Cordonnier, and Andreas Loukas. Attention is not all you need: Pure attention loses rank doubly exponentially with depth. In *International Conference on Machine Learning*, pp. 2793–2803. PMLR, 2021.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2020.

- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of machine learning research*, 12(7), 2011.
- Aaron Gokaslan and Vanya Cohen. Openwebtext corpus.
- Gene H Golub and Charles F Van Loan. *Matrix computations*. JHU press, 2013.
- Alexander Graham. *Kronecker products and matrix calculus with applications*. Courier Dover Publications, 2018.
- Dongchen Han, Xuran Pan, Yizeng Han, Shiji Song, and Gao Huang. Flatten transformer: Vision transformer using focused linear attention. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 5961–5971, 2023.
- Alex Henry, Prudhvi Raj Dachapally, Shubham Shantaram Pawar, and Yuxuan Chen. Query-key normalization for transformers. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pp. 4246–4253, 2020.
- Roger A Horn and Charles R Johnson. Topics in matrix analysis, 1991. *Cambridge University Press, Cambridge*, 37:39, 1991.
- Roger A Horn and Charles R Johnson. *Matrix analysis*. Cambridge university press, 2012.
- Andrej Karpathy. NanoGPT. <https://github.com/karpathy/nanoGPT>, 2022.
- Hyunjik Kim, George Papamakarios, and Andriy Mnih. The lipschitz constant of self-attention. In *International Conference on Machine Learning*, pp. 5562–5571. PMLR, 2021.
- Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Atli Kosson, Bettina Messmer, and Martin Jaggi. Rotational equilibrium: How weight decay balances learning across neural networks. *arXiv preprint arXiv:2305.17212*, 2023.
- Atli Kosson, Bettina Messmer, and Martin Jaggi. Analyzing & reducing the need for learning rate warmup in gpt training. *Advances in Neural Information Processing Systems*, 2024.
- Yann LeCun, Bernhard Boser, John S Denker, Donnie Henderson, Richard E Howard, Wayne Hubbard, and Lawrence D Jackel. Backpropagation applied to handwritten zip code recognition. *Neural computation*, 1(4):541–551, 1989.
- Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- Yann LeCun, Léon Bottou, Genevieve B Orr, and Klaus-Robert Müller. Efficient backprop. In *Neural networks: Tricks of the trade*, pp. 9–50. Springer, 2002.
- Liyuan Liu, Xiaodong Liu, Jianfeng Gao, Weizhu Chen, and Jiawei Han. Understanding the difficulty of training transformers. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 5747–5763, 2020.
- Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 10012–10022, 2021.
- Ilya Loshchilov and Frank Hutter. Sgdr: Stochastic gradient descent with warm restarts. *arXiv preprint arXiv:1608.03983*, 2016.
- Ilya Loshchilov and Frank Hutter. Fixing weight decay regularization in adam. In *International Conference on Learning Representations*, 2019.

- Lorenzo Noci, Sotiris Anagnostidis, Luca Biggio, Antonio Orvieto, Sidak Pal Singh, and Aurelien Lucchi. Signal propagation in transformers: Theoretical perspectives and the role of rank collapse. *Advances in Neural Information Processing Systems*, 35:27198–27211, 2022.
- William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 4195–4205, 2023.
- Kaare Brandt Petersen, Michael Syskind Pedersen, et al. The matrix cookbook. *Technical University of Denmark*, 7(15):510, 2008.
- Xianbiao Qi, Jianan Wang, Yihao Chen, Yukai Shi, and Lei Zhang. Lipsformer: Introducing lipschitz continuity to vision transformers. In *The Eleventh International Conference on Learning Representations*, 2023a.
- Xianbiao Qi, Jianan Wang, and Lei Zhang. Understanding optimization of deep learning via jacobian matrix and lipschitz constant. *arXiv preprint arXiv:2306.09338*, 2023b.
- Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. Improving language understanding by generative pre-training. 2018.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pp. 8748–8763. PMLR, 2021.
- Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *International conference on machine learning*, pp. 8821–8831. Pmlr, 2021.
- Herbert Robbins and Sutton Monro. A stochastic approximation method. *The annals of mathematical statistics*, pp. 400–407, 1951.
- David E Rumelhart, Geoffrey E Hinton, and Ronald J Williams. Learning representations by back-propagating errors. *nature*, 323(6088):533–536, 1986.
- Leslie N Smith. Cyclical learning rates for training neural networks. In *2017 IEEE winter conference on applications of computer vision (WACV)*, pp. 464–472. IEEE, 2017.
- Yuandong Tian, Yiping Wang, Beidi Chen, and Simon S Du. Scan and snap: Understanding training dynamics and token composition in 1-layer transformer. *Advances in Neural Information Processing Systems*, 36:71911–71947, 2023a.
- Yuandong Tian, Yiping Wang, Zhenyu Zhang, Beidi Chen, and Simon Du. Joma: Demystifying multilayer transformers via joint dynamics of mlp and attention. *arXiv preprint arXiv:2310.00535*, 2023b.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- Roman Vershynin. *High-dimensional probability: An introduction with applications in data science*, volume 47. Cambridge university press, 2018.
- Hongyu Wang, Shuming Ma, Li Dong, Shaohan Huang, Dongdong Zhang, and Furu Wei. Deepnet: Scaling transformers to 1,000 layers. *arXiv preprint arXiv:2203.00555*, 2022.

- Qiang Wang, Bei Li, Tong Xiao, Jingbo Zhu, Changliang Li, Derek F Wong, and Lidia S Chao. Learning deep transformer models for machine translation. *arXiv preprint arXiv:1906.01787*, 2019.
- Ross Wightman. Pytorch image models. <https://github.com/rwightman/pytorch-image-models>, 2019.
- John Wright and Yi Ma. *High-dimensional data analysis with low-dimensional models: Principles, computation, and applications*. Cambridge University Press, 2022.
- Ruibin Xiong, Yunchang Yang, Di He, Kai Zheng, Shuxin Zheng, Chen Xing, Huishuai Zhang, Yanyan Lan, Liwei Wang, and Tieyan Liu. On layer normalization in the transformer architecture. In *International Conference on Machine Learning*, pp. 10524–10533. PMLR, 2020.
- Shuangfei Zhai, Tatiana Likhomanenko, Etai Littwin, Dan Busbridge, Jason Ramapuram, Yizhe Zhang, Jiatao Gu, and Joshua M Susskind. Stabilizing transformer training by preventing attention entropy collapse. In *International Conference on Machine Learning*, pp. 40770–40803. PMLR, 2023.
- Biao Zhang and Rico Sennrich. Root mean square layer normalization. *Advances in Neural Information Processing Systems*, 32, 2019.

A KRONECKER PRODUCT AND VECTORIZATION

Kronecker product (Graham, 2018; Petersen et al., 2008), also called as matrix direct product, is an operation defined on two matrices of arbitrary size. The specific definition is as follows.

Definition 1 (Kronecker Product)

Let A be an $n \times p$ matrix and B an $m \times q$ matrix. The $mn \times pq$ matrix

$$A \otimes B = \begin{bmatrix} a_{1,1}B & a_{1,2}B & \cdots & a_{1,p}B \\ a_{2,1}B & a_{2,2}B & \cdots & a_{2,p}B \\ \vdots & \vdots & \ddots & \vdots \\ a_{n,1}B & a_{n,2}B & \cdots & a_{n,p}B \end{bmatrix}$$

is called the Kronecker product of A and B . It is also called the direct product or the tensor product.

For instance, if $A = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}$, and $B = \begin{bmatrix} 1 & 2 & 3 \\ 3 & 4 & 5 \end{bmatrix}$, then $A \otimes B = \begin{bmatrix} 1 & 2 & 3 & 2 & 4 & 6 \\ 3 & 4 & 5 & 6 & 8 & 10 \\ 3 & 6 & 9 & 4 & 8 & 12 \\ 9 & 12 & 15 & 12 & 16 & 20 \end{bmatrix}$.

Some basic properties of the Kronecker product include:

$$\begin{aligned} A \otimes (B \otimes C) &= (A \otimes B) \otimes C, \\ A \otimes (B + C) &= (A \otimes B) + (A \otimes C), \quad (A + B) \otimes C = (A \otimes C) + (B \otimes C), \\ (A \otimes B)^T &= A^T \otimes B^T. \end{aligned}$$

For a matrix A , the rank of $A \otimes A$ can be computed as,

$$\text{rank}(A \otimes A) = \text{rank}(A) \cdot \text{rank}(A).$$

It means if the rank of the matrix A is small, then the rank of $A \otimes A$ will also be very small.

In mathematics, *Vectorization* (Graham, 2018; Petersen et al., 2008) is usually used together with the Kronecker product to express matrix multiplication as a linear transformation on matrices. After vectorization, we can calculate the Jacobian matrix of the matrix product more conveniently. A property of vectorization for the matrix product is defined below.

Proposition 2 (Property of Vectorization for Matrix Product)

Let $A \in \mathcal{R}^{m \times n}$, $B \in \mathcal{R}^{n \times k}$, $C \in \mathcal{R}^{k \times l}$, then we have

$$\text{vec}(ABC) = (C^T \otimes A) \text{vec}(B).$$

Proof. Let C_i be the i -th row of C . Then we have:

$$\begin{aligned} \text{vec}(ABC) &= \sum_{i=1}^n \sum_{j=1}^k b_{ij} \text{vec}(a_i C_j) \\ &= \sum_{i=1}^n \sum_{j=1}^k b_{ij} (C_j^T \otimes a_i) \\ &= \sum_{j=1}^k (C_j^T \otimes A) b_j \\ &= (C^T \otimes A) \text{vec}(B). \end{aligned}$$

□

Furthermore, we have the following properties:

$$\begin{aligned} \text{vec}(AB) &= (I_k \otimes A) \text{vec}(B) = (B^T \otimes I_m) \text{vec}(A), \\ \text{vec}(ABC) &= (C^T B^T \otimes I_m) \text{vec}(A) \\ &= (C^T \otimes A) \text{vec}(B) \\ &= (I_l \otimes AB) \text{vec}(C) \end{aligned}$$

where $\mathbf{I}_k \in \mathcal{R}^{k \times k}$, $\mathbf{I}_l \in \mathcal{R}^{l \times l}$, $\mathbf{I}_m \in \mathcal{R}^{m \times m}$ are all identity matrices.

Together with the Kronecker product, vectorization is an effective tool to compute matrix calculus. We can see the following two examples.

Let $\mathbf{P} = \mathbf{AB}$ where $\mathbf{A} \in \mathcal{R}^{m \times n}$, $\mathbf{B} \in \mathcal{R}^{n \times k}$, we have:

$$\frac{\partial \text{vec}(\mathbf{P})}{\partial \text{vec}(\mathbf{A})} = \mathbf{B}^\top \otimes \mathbf{I}_m, \quad \frac{\partial \text{vec}(\mathbf{P})}{\partial \text{vec}(\mathbf{B})} = \mathbf{I}_k \otimes \mathbf{A}.$$

Let $\mathbf{P} = \mathbf{ABC}$ where $\mathbf{A} \in \mathcal{R}^{m \times n}$, $\mathbf{B} \in \mathcal{R}^{n \times k}$, $\mathbf{C} \in \mathcal{R}^{k \times l}$, we have,

$$\frac{\partial \text{vec}(\mathbf{P})}{\partial \text{vec}(\mathbf{A})} = \mathbf{C}^\top \mathbf{B}^\top \otimes \mathbf{I}_m, \quad \frac{\partial \text{vec}(\mathbf{P})}{\partial \text{vec}(\mathbf{B})} = \mathbf{C}^\top \otimes \mathbf{A}, \quad \frac{\partial \text{vec}(\mathbf{P})}{\partial \text{vec}(\mathbf{C})} = \mathbf{I}_l \otimes \mathbf{AB}.$$

Vectorization and Kronecker product provide us a convenient way to analyze the self-attention module. We can compute the Jacobian matrix of the output with respect to the input or the weight matrix more conveniently. For more introduction to Kronecker product and vectorization, the readers can refer to (Petersen et al., 2008; Graham, 2018)

B DERIVATION OF JACOBIAN MATRIX FOR SINGLE-HEAD SELF-ATTENTION

A single-head self-attention can be defined as

$$\mathbf{Y} = \mathbf{W}_v \mathbf{X} \mathbf{A},$$

where $\mathbf{A} = \text{softmax}(\frac{\mathbf{P}}{\sqrt{d_q}})$ in which $\mathbf{P} = \mathbf{X}^\top \mathbf{W}_q^\top \mathbf{W}_k \mathbf{X}$. The matrix \mathbf{A} is called the attention matrix and $\frac{\mathbf{P}}{\sqrt{d_q}}$ is called the logit, $\mathbf{A} \in \mathcal{R}^{n \times n}$, $\mathbf{X} \in \mathcal{R}^{d \times n}$, $\mathbf{W}_v \in \mathcal{R}^{d_v \times d}$. Here, our goal is to calculate $\frac{\partial \text{vec}(\mathbf{Y})}{\partial \text{vec}(\mathbf{X})}$.

In the main body, we have derived $\frac{\partial \text{vec}(\mathbf{P})}{\partial \text{vec}(\mathbf{X})}$. Here, let us calculate the matrix calculus of $\mathbf{A} = \text{softmax}(\frac{\mathbf{P}}{\sqrt{d_q}})$ with respect to \mathbf{P} using Kronecker products and vectorization. We can rewrite it as $\mathbf{A} = \exp(\frac{\mathbf{P}}{\sqrt{d_q}}) \oslash (\mathbf{1}_n \mathbf{1}_n^\top \exp(\frac{\mathbf{P}}{\sqrt{d_q}}))$, where $\mathbf{1}_n$ denotes an n -dimensional vector of 1's in \mathcal{R}^n . Note that \mathbf{A} is obtained by conducting a softmax operation in each column individually.

First, let us define two intermediate variables:

$$\mathbf{B} = \exp(\frac{\mathbf{P}}{\sqrt{d_q}}), \quad \mathbf{C} = \mathbf{1} \mathbf{1}^\top \exp(\frac{\mathbf{P}}{\sqrt{d_q}}) = \mathbf{1} \mathbf{1}^\top \mathbf{B}.$$

In this way, we can represent the attention matrix \mathbf{A} as $\mathbf{A} = \mathbf{B} \oslash \mathbf{C}$.

Then, we can vectorize the equation $\mathbf{A} = \mathbf{B} \oslash \mathbf{C}$ as follows:

$$\text{vec}(\mathbf{A}) = \text{vec}(\mathbf{B} \oslash \mathbf{C}) = \text{vec}(\mathbf{B}) \oslash \text{vec}(\mathbf{C}).$$

According to the chain rule, we have

$$\frac{\partial \text{vec}(\mathbf{A})}{\partial \text{vec}(\mathbf{P})} = \frac{\partial \text{vec}(\mathbf{A})}{\partial \text{vec}(\mathbf{B})} \frac{\partial \text{vec}(\mathbf{B})}{\partial \text{vec}(\mathbf{P})} + \frac{\partial \text{vec}(\mathbf{A})}{\partial \text{vec}(\mathbf{C})} \frac{\partial \text{vec}(\mathbf{C})}{\partial \text{vec}(\mathbf{P})}.$$

Let us calculate each individual term. We have

$$\begin{aligned}\frac{\partial \text{vec}(\mathbf{A})}{\partial \text{vec}(\mathbf{B})} &= \mathbf{1}_{n^2} \otimes \text{diag}(\text{vec}(\mathbf{C})), \\ \frac{\partial \text{vec}(\mathbf{B})}{\partial \text{vec}(\mathbf{P})} &= \frac{\text{diag}(\text{vec}(\mathbf{B}))}{\sqrt{d_q}}, \\ \frac{\partial \text{vec}(\mathbf{A})}{\partial \text{vec}(\mathbf{C})} &= -\text{diag}(\text{vec}(\mathbf{B}) \otimes (\text{vec}(\mathbf{C}) \odot \text{vec}(\mathbf{C}))), \\ \frac{\partial \text{vec}(\mathbf{C})}{\partial \text{vec}(\mathbf{P})} &= \frac{(\mathbf{I}_n \otimes \mathbf{1}_n \mathbf{1}_n^\top) \text{diag}(\text{vec}(\mathbf{B}))}{\sqrt{d_q}},\end{aligned}$$

where \mathbf{I}_n is an identity matrix of $n \times n$ and \otimes is the Kronecker product, $\mathbf{1}_{nn}$ denotes an n^2 -dimensional vector of 1's in \mathcal{R}^{n^2} .

By substituting the above four terms into the chain rule, we have

$$\begin{aligned}\frac{\partial \text{vec}(\mathbf{A})}{\partial \text{vec}(\mathbf{P})} &= \frac{(\mathbf{1}_{n^2} \otimes \text{diag}(\text{vec}(\mathbf{C}))) \text{diag}(\text{vec}(\mathbf{B})) - \text{diag}(\text{vec}(\mathbf{B}) \otimes (\text{vec}(\mathbf{C}) \odot \text{vec}(\mathbf{C}))) (\mathbf{I}_n \otimes \mathbf{1}_n \mathbf{1}_n^\top) \text{diag}(\text{vec}(\mathbf{B}))}{\sqrt{d_q}} \\ &= \frac{\text{diag}(\text{vec}(\mathbf{A})) - \text{diag}(\text{vec}(\mathbf{B}) \otimes (\text{vec}(\mathbf{C}) \odot \text{vec}(\mathbf{C}))) (\mathbf{I}_n \otimes \mathbf{1}_n \mathbf{1}_n^\top) \text{diag}(\text{vec}(\mathbf{B}))}{\sqrt{d_q}} \\ &= \frac{\text{blockdiag}(\text{diag}(\mathbf{A}_{:,1}) - \mathbf{A}_{:,1} \mathbf{A}_{:,1}^\top, \dots, \text{diag}(\mathbf{A}_{:,n}) - \mathbf{A}_{:,n} \mathbf{A}_{:,n}^\top)}{\sqrt{d_q}}.\end{aligned}$$

For the simplicity, we denote

$$\mathbf{J} = \text{blockdiag}(\text{diag}(\mathbf{A}_{:,1}) - \mathbf{A}_{:,1} \mathbf{A}_{:,1}^\top, \dots, \text{diag}(\mathbf{A}_{:,n}) - \mathbf{A}_{:,n} \mathbf{A}_{:,n}^\top).$$

When \mathbf{A} and \mathbf{P} are vectorized into vectors, we use \mathbf{a} and \mathbf{p} to denote them, respectively. Then we see that

$$\frac{\partial \text{vec}(\mathbf{a})}{\partial \text{vec}(\mathbf{p})} = \frac{\text{diag}(\mathbf{a}) - \mathbf{a} \mathbf{a}^\top}{\sqrt{d_q}}.$$

If \mathbf{a} approaches to a unit vector \mathbf{e} , then the Jacobian matrix $\frac{\partial \text{vec}(\mathbf{a})}{\partial \text{vec}(\mathbf{p})}$ will tend to $\mathbf{0}$.

In Section 3.2, we have the following Jacobian matrix

$$\frac{\partial \text{vec}(\mathbf{P})}{\partial \text{vec}(\mathbf{X})} = (\mathbf{X}^\top \mathbf{W}_k^\top \mathbf{W}_q \otimes \mathbf{I}_n) \mathbf{K} + (\mathbf{I}_n \otimes \mathbf{X}^\top \mathbf{W}_q^\top \mathbf{W}_k).$$

By vectorization of $\mathbf{Y} = \mathbf{W}_v \mathbf{X} \mathbf{A}$, we have

$$\partial \text{vec}(\mathbf{Y}) = (\mathbf{A}^\top \otimes \mathbf{W}_v) \partial \text{vec}(\mathbf{X}) + (\mathbf{I}_n \otimes \mathbf{W}_v \mathbf{X}) \partial \text{vec}(\mathbf{A}).$$

Therefore, according to the product rule and chain rule, we can denote the Jacobian matrix of \mathbf{Y} with respect to \mathbf{X} as follows:

$$\begin{aligned}\frac{\partial \text{vec}(\mathbf{Y})}{\partial \text{vec}(\mathbf{X})} &= (\mathbf{A}^\top \otimes \mathbf{W}_v) + (\mathbf{I}_n \otimes \mathbf{W}_v \mathbf{X}) \frac{\partial \text{vec}(\mathbf{A})}{\partial \text{vec}(\mathbf{X})}, \\ &= (\mathbf{A}^\top \otimes \mathbf{W}_v) + (\mathbf{I}_n \otimes \mathbf{W}_v \mathbf{X}) \frac{\partial \text{vec}(\mathbf{A})}{\partial \text{vec}(\mathbf{P})} \frac{\partial \text{vec}(\mathbf{P})}{\partial \text{vec}(\mathbf{X})}.\end{aligned}$$

Bringing in all these terms, we get the following formula:

$$\frac{\partial \text{vec}(\mathbf{Y})}{\partial \text{vec}(\mathbf{X})} = (\mathbf{A}^\top \otimes \mathbf{W}_v) + (\mathbf{I}_n \otimes \mathbf{W}_v \mathbf{X}) \frac{\mathbf{J}}{\sqrt{d_q}} \left((\mathbf{X}^\top \mathbf{W}_k^\top \mathbf{W}_q \otimes \mathbf{I}_n) \mathbf{K} + (\mathbf{I}_n \otimes \mathbf{X}^\top \mathbf{W}_q^\top \mathbf{W}_k) \right). \quad (7)$$

Let us analyze Equation 7. If a malignant entropy mode happens, $\frac{J}{\sqrt{d_q}}$ will approach 0 because each \mathbf{a} in \mathbf{A} will be a unit vector \mathbf{e} . From the perspective of the forward process, the features \mathbf{Y} will collapse to several directions. From the perspective of the backward process, $\frac{J}{\sqrt{d_q}}$ will become 0, and $\frac{\partial \text{vec}(\mathbf{Y})}{\partial \text{vec}(\mathbf{X})}$ will be a sparse and simultaneously low-rank matrix. Through $\mathbf{A}^\top \otimes \mathbf{W}_v$, most of the positions in \mathbf{X} will get zero gradient, and only very few columns will obtain some large noisy gradients. In a malignant entropy mode, the learned feature is invalid and useless. Similarly, if a benign entropy mode happens, the attention map \mathbf{A} will approach an identity matrix \mathbf{I} and $\frac{J}{\sqrt{d_q}} \approx 0$ when $\mathbf{A} \approx \mathbf{I}$. Therefore, we have $\frac{\partial \text{vec}(\mathbf{Y})}{\partial \text{vec}(\mathbf{X})} \approx (\mathbf{I}^\top \otimes \mathbf{W}_v)$. In this way, a self-attention module degenerates to a linear layer.

C PROOF OF BENIGN ENTROPY COLLAPSE

Recall that $\mathbf{A} = \text{softmax}(\frac{\mathbf{P}}{\sqrt{d_q}})$ where $\mathbf{P} = \mathbf{X}^\top \mathbf{W}_q^\top \mathbf{W}_k \mathbf{X}$. Here, let $\mathbf{W} = \mathbf{W}_q^\top \mathbf{W}_k$ and $\mathbf{W} \in \mathcal{R}^{d \times d}$. In this way, we have that $\mathbf{A} = \text{softmax}(\frac{\mathbf{X}^\top \mathbf{W} \mathbf{X}}{\sqrt{d_q}})$. We know that $\text{rank}(\mathbf{W}) \leq d_q$ and $d_q < d$. To prove \mathbf{A} will always collapse to an identity matrix when \mathbf{W} is a non-symmetric positive quasi-definite square matrix, it is equivalent to prove $\mathbb{E}[\mathbf{x}_i^\top \mathbf{W} \mathbf{x}_i] \gg \mathbb{E}[\mathbf{x}_i^\top \mathbf{W} \mathbf{x}_j]$ for any $i \neq j$. It will be very hard to prove it mathematically if \mathbf{W} is a form of a non-symmetric positive quasi-definite square matrix. Therefore, let us make some simplification assumptions. Assume \mathbf{W} is a real symmetric positive semi-definite square matrix and its trace is in direct proportion to the dimension d_q , and any \mathbf{x}_i is a high-dimension random vector and each element in $\mathbf{x}_{i,j} \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1)$.

We break our proof into two sub-problems.

Proposition 3 (Expectation of $\mathbf{x}_i^\top \mathbf{W} \mathbf{x}_i$)

Let \mathbf{W} be a real symmetric positive semi-definite matrix, and any \mathbf{x}_i be a high-dimensional random vector. Then, we have $\mathbb{E}[\mathbf{x}_i^\top \mathbf{W} \mathbf{x}_i] = \text{trace}(\mathbf{W})$.

Proof. Let \mathbf{W} be a real symmetric positive semi-definite, thus it can be decomposed into $\mathbf{W} = \mathbf{U} \Sigma \mathbf{U}^\top$. In this way, we have

$$\begin{aligned} \mathbb{E}[\mathbf{x}_i^\top \mathbf{W} \mathbf{x}_i] &= \mathbb{E}[\mathbf{x}_i^\top \mathbf{U} \Sigma \mathbf{U}^\top \mathbf{x}_i] \\ &= \mathbb{E}[\mathbf{z}^\top \Sigma \mathbf{z}] \quad (\text{let } \mathbf{z} = \mathbf{U}^\top \mathbf{x}_i) \\ &= \mathbb{E}\left[\sum_{i=1}^d \sigma_i z_i^2\right] \quad (\Sigma \text{ is a diagonal matrix}) \\ &= \sum_{i=1}^{d_q} \sigma_i \times (0 + 1) \quad (\text{by independence, mean 0}) \\ &= \sum_{i=1}^{d_q} \sigma_i = \text{trace}(\mathbf{W}). \end{aligned}$$

For a real symmetric positive semi-definite, all its singular values are larger or equal to 0. Thus, we know $\sum_{i=1}^{d_q} \sigma_i > 0$ considering \mathbf{W} is not a matrix of all zeros. \square

Proposition 4 (Expectation of $\mathbf{x}_i^\top \mathbf{W} \mathbf{x}_j$ for $i \neq j$)

Let \mathbf{W} be a real symmetric positive semi-definite square matrix, and any \mathbf{x}_i is a high-dimension random vector; $\mathbb{E}_{i \neq j}[\mathbf{x}_i^\top \mathbf{W} \mathbf{x}_j] = 0$.

Proof. Let \mathbf{W} be a real symmetric positive semi-definite. Thus, it can be decomposed into $\mathbf{W} = \mathbf{U}\Sigma\mathbf{U}^\top$. In this way, we have

$$\begin{aligned}\mathbb{E}_{i \neq j} [\mathbf{x}_i^\top \mathbf{W} \mathbf{x}_j] &= \mathbb{E}_{i \neq j} [\mathbf{x}_i^\top \mathbf{U} \Sigma \mathbf{U}^\top \mathbf{x}_j] \\ &= \mathbb{E} [\mathbf{z}^\top \Sigma \mathbf{v}] \quad (\text{let } \mathbf{z} = \mathbf{U}^\top \mathbf{x}_i, \text{ and } \mathbf{v} = \mathbf{U}^\top \mathbf{x}_j) \\ &= \mathbb{E} \left[\sum_{i=1}^{d_q} \sigma_{ij} z_i v_j \right] \quad (\Sigma \text{ is a diagonal matrix. } \mathbf{z} \text{ and } \mathbf{v} \text{ are independent}) \\ &= 0.\end{aligned}$$

□

According to Proposition 3 and Proposition 4, we can have that $\mathbb{E} [\mathbf{x}_i^\top \mathbf{W} \mathbf{x}_i] > \mathbb{E} [\mathbf{x}_i^\top \mathbf{W} \mathbf{x}_j]$ for any $i \neq j$. Considering that \mathbf{W} is usually a high-dimensional matrix and some of its singular values are significantly larger than 0. Thus, after the softmax operation, \mathbf{A} will always collapse to an identity matrix. In this way, the self-attention module degenerates into a linear projection module. The model fitting ability will decline, but model training will not crash. Our proof is based on matrix computations (Golub & Van Loan, 2013) and high-dimensional probability (Vershynin, 2018).

D PROOF OF MALIGNANT ENTROPY COLLAPSE

Proof. Step 1. To prove the sparsity of \mathbf{A} , we must show that the number of non-zero elements in each column of \mathbf{A} is small.

Now, let's consider the properties of the matrix $\mathbf{P} = \mathbf{X}^\top \mathbf{W} \mathbf{X}$. Since \mathbf{X} and \mathbf{W} are low-rank matrices, the matrix \mathbf{P} will also be low-rank. Specifically, the rank of \mathbf{P} is bounded by the rank of \mathbf{X} and \mathbf{W} , which is much smaller than the dimension. In particular, \mathbf{P} has only a small number of significant singular values. This implies that the entries in \mathbf{P} are concentrated in a lower-dimensional subspace.

When we apply the softmax function to the rows of \mathbf{P} , the function concentrates most of the probability mass on a few components of each column. This is because the softmax function is sharply peaked around the largest values in each row. The smaller values in each row contribute less to the sum in the denominator of the softmax function, and therefore, their corresponding entries in \mathbf{A} will be small.

Thus, in each column of \mathbf{A} , only a small number of entries will be non-zero with high probability, and the rest will be close to zero. This establishes the sparsity of \mathbf{A} .

Step 2. To prove the low-rankness of \mathbf{A} , we turn to proving that \mathbf{A} is approximately low-rank.

As noted earlier, the matrix $\mathbf{P} = \mathbf{X}^\top \mathbf{W} \mathbf{X}$ is low-rank. Specifically, the rank of \mathbf{P} is bounded by the ranks of \mathbf{X} and \mathbf{W} , which are both small. Therefore, \mathbf{P} has only a few dominant singular values. Assume that the softmax function does not significantly change the rank of the matrix. The rank of \mathbf{A} is controlled by the rank of \mathbf{P} , as the softmax operation only introduces nonlinearities that do not increase the rank.

Thus, since \mathbf{P} has a small number of dominant singular values, \mathbf{A} , formed by the softmax of \mathbf{P} , will also have a small number of significant singular values. This implies that \mathbf{A} is approximately low-rank.

In conclusion, under the assumptions that \mathbf{X} and \mathbf{W} are low-rank matrices with sufficiently large singular values, the matrix \mathbf{A} formed by applying the softmax function to \mathbf{P} will exhibit the properties of both sparsity and low-rankness with high probability.

□

E PROOF OF WEYL’S INEQUALITY ON SINGULAR VALUES

Our derivation depends on Horn & Johnson (1991; 2012). Readers can refer to these material for more background information.

Proof. Before we prove Weyl’s Inequality on singular values, let us review Courant-Fischer min-max principle that is important for analyzing the singular values of matrix.

Theorem 3 (Courant-Fischer Min-max Principle for Singular Values)

Let $\mathbf{W} \in \mathcal{R}^{m \times n}$ be a matrix where $m \geq n$. \mathbf{W} has ordered singular values $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n \geq 0$. Then, for $i = 1, 2, \dots, n$, we have

$$\sigma_i(\mathbf{W}) = \max_{\substack{S \subset \mathcal{R}^n \\ \dim(S)=i}} \min_{\substack{\mathbf{x} \in S \\ \|\mathbf{x}\|=1}} \|\mathbf{W}\mathbf{x}\| = \min_{\substack{S' \subset \mathcal{R}^n \\ \dim(S')=n-i+1}} \max_{\substack{\mathbf{x} \in S' \\ \|\mathbf{x}\|=1}} \|\mathbf{W}\mathbf{x}\|$$

where the maximum is taken over all i -dimensional subspaces S of \mathcal{R}^n , and the minimum is taken over all unit vectors \mathbf{x} in S .

Let $\mathbf{W}_1 = \mathbf{U}\mathbf{\Sigma}_1\mathbf{V}^\top$ and $\mathbf{W}_2 = \mathbf{U}\mathbf{\Sigma}_2\mathbf{V}^\top$ be singular value decompositions of \mathbf{W}_1 and \mathbf{W}_2 with unitary matrix $\mathbf{V} = [\mathbf{v}_1, \dots, \mathbf{v}_n]$, $\mathbf{V} = [\mathbf{v}_1, \dots, \mathbf{v}_n]$ where $\mathbf{v}_i, \mathbf{v}_i \in \mathcal{R}^n$ and unitary matrix $\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_m]$, $\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_m]$, where $\mathbf{u}_j, \mathbf{u}_j \in \mathcal{R}^m$.

Let i and j be positive integers with $1 \leq i, j \leq n$ and $i + j \leq n + 1$. Let $S_1 \equiv \text{Span}\{\mathbf{v}_i, \dots, \mathbf{v}_n\}$ and $S_2 \equiv \text{Span}\{\mathbf{v}_j, \dots, \mathbf{v}_n\}$; notice that $\dim(S_1) = n - i + 1$ and $\dim(S_2) = n - j + 1$. Let $k \equiv \dim(S_1 \cap S_2)$, then we have

$$\begin{aligned} \dim(S_1 \cap S_2) &= \dim(S_1) + \dim(S_2) - \dim(S_1 + S_2) = (n - i + 1) + (n - j + 1) - \dim(S_1 + S_2) \\ &\geq (n - i + 1) + (n - j + 1) - n = n - (i + j - 1) + 1 \geq 1. \end{aligned}$$

Because of the bounds assumed for i and j . Thus, the subspace $S_1 \cap S_2$ has positive dimension k , $n - k + 1 \leq i + j - 1$, and we have

$$\begin{aligned} \sigma_{i+j-1}(\mathbf{W}_1 + \mathbf{W}_2) &\leq \sigma_{n-k+1}(\mathbf{W}_1 + \mathbf{W}_2) \\ &= \min_{\substack{S \subset \mathcal{R}^n \\ \dim(S)=k}} \max_{\substack{\mathbf{x} \in S \\ \|\mathbf{x}\|_2=1}} \|(\mathbf{W}_1 + \mathbf{W}_2)\mathbf{x}\|_2 \\ &\leq \max_{\substack{\mathbf{x} \in S_1 \cap S_2 \\ \|\mathbf{x}\|_2=1}} \|(\mathbf{W}_1 + \mathbf{W}_2)\mathbf{x}\|_2 \\ &\leq \max_{\substack{\mathbf{x} \in S_1 \cap S_2 \\ \|\mathbf{x}\|_2=1}} \|\mathbf{W}_1\mathbf{x}\|_2 + \max_{\substack{\mathbf{x} \in S_1 \cap S_2 \\ \|\mathbf{x}\|_2=1}} \|\mathbf{W}_2\mathbf{x}\|_2 \\ &\leq \max_{\substack{\mathbf{x} \in S_1 \\ \|\mathbf{x}\|_2=1}} \|\mathbf{W}_1\mathbf{x}\|_2 + \max_{\substack{\mathbf{x} \in S_2 \\ \|\mathbf{x}\|_2=1}} \|\mathbf{W}_2\mathbf{x}\|_2 = \sigma_i(\mathbf{W}_1) + \sigma_j(\mathbf{W}_2). \end{aligned}$$

As a special case, the second part of the theorem follows directly from the general result of part (a). Specifically, for $i = j = 1$, we have:

$$\sigma_1(\mathbf{W}_1 + \mathbf{W}_2) \leq \sigma_1(\mathbf{W}_1) + \sigma_1(\mathbf{W}_2).$$

This completes the proof. \square

F RELATED WORKS

Training Dynamics of Transformer. Previous works have delved into understanding the training dynamics of Transformers from two different perspectives: a high-level perspective and a low-level perspective. From a high-level perspective, Scan&Snap (Tian et al., 2023a) unveiled complex phenomena, particularly in single-layer architectures, relating to frequency and discriminative bias. These studies linked sparse attention patterns to token co-occurrence frequencies and observed two-stage behaviors in attention logits. JoMA (Tian et al., 2023b) further improved upon previous models by incorporating residual connections and MLP nonlinearity, analyzing joint training of MLP and self-attention layers, and offering qualitative explanations for multi-layer Transformer dynamics. From a low-level perspective, two critical challenges in Transformer training have been identified: rank collapse (Dong et al., 2021; Noci et al., 2022), where attention output converges to a rank 1 matrix, potentially causing vanishing gradients; and entropy collapse (Zhai et al., 2023), which denotes pathologically low attention entropy, corresponding to highly concentrated attention scores. *In this work, we analyze and prove two different entropy collapse modes and identify the key reason for model failure is spectral energy concentration. Finally, we introduce a simple but effective solution to address this problem.*

Training Stability of Transformer. ReZero (Bachlechner et al., 2021) introduces a simple yet effective mechanism for improving training stability. The key innovation lies in initializing residual connections to zero, which allows networks to learn identity mappings more easily. Admin (Liu et al., 2020) introduces a new network initialization strategy tailored for Transformers to make the network train stable. DeepNorm (Wang et al., 2022) extends the concept of normalization to accommodate increasingly deeper networks. By dynamically adjusting normalization parameters, DeepNorm ensures stability even as network depth increases. LipsFormer Qi et al. (2023a) addresses the specific challenge of stability in transformer networks. By introducing a Lipschitz continuity constraint, Lipsformer effectively mitigates the issue of exploding gradients - a common problem in deep transformer architectures. This approach ensures that the network's output changes smoothly with respect to its input, promoting overall stability. ReZero, Admin, and DeepNorm can all be considered as an approach to control the Lipschitz constant of the network in the initial stage. *In this work, by revisiting the training dynamics of Transformer, we can achieve a stable training only by modifying the optimizer instead of using learning rate warmup or changing the network structures as LipsFormer (Qi et al., 2023a) and QKNorm Henry et al. (2020); Dehghani et al. (2023).*

Learning Rate Schedule. Warmup (Loshchilov & Hutter, 2016) has emerged as a must-have technique for ensuring a stable network training, especially in the initial phases of the optimization process. This method involves gradually increasing the learning rate from a small value to the desired initial learning rate over a certain number of training steps or epochs. The cosine learning rate scheduler (Loshchilov & Hutter, 2016) has gained popularity due to its smooth annealing properties. This schedule decreases the learning rate following a cosine curve, starting from an initial value and decaying to a minimum value over a set number of epochs or iterations. Cyclic learning rates (Smith, 2017) involve systematically varying the learning rate between boundary values. The learning rate oscillates between a lower and upper bound, either linearly or following other patterns (e.g., triangular, cosine). The above-mentioned learning rate schedules require specification of a stopping time step T , Defazio et al. (2024) introduces a Schedule-Free approach that avoids the need for this stopping time by eschewing the use of schedules entirely.

Compared to the up-mentioned works, the core novel contributions of our work lie on follows.

1. We present a theoretical analysis for Transformer training and point out two entropy collapse modes, *i.e.* the benign collapse and the malignant collapse.
2. We reveal that *spectral energy concentration (SEC)* of $\mathbf{W}_q^\top \mathbf{W}_k$ is the main reason of model crash.
3. We introduce AdamW², a new optimization strategy motivated by Weyl's Inequality.

We also observe there are two works (Kosson et al., 2023; 2024) discussing the needs of learning rate warmup by explicitly controlling the angular updates via Rotational Optimizer Variants and by limiting the Frobenius norm of the update relative to that of the weights. They provided some different perspectives on the necessity of the learning rate warmup.

G SIMULATION OF THREE ATTENTION MODES

We provide a simple simulation code to simulate three attention modes, but it is important to note that the real picture is more complicated. In real case, in the benign attention entropy mode, $W_q^\top W$ is a non-symmetric positive quasi-definite square matrix instead of a symmetric positive definite matrix in our simulation. The code is just to demonstrate the core ideas behind three attention modes.

CODE 1: Simulation of Three Attention Modes.

```

1 import torch
2 import torch.nn
3 import matplotlib.pyplot as plt
4 import numpy as np
5
6
7 #Randomly generate data and weight matrices
8 d_q, d, num_tokens= 64, 768, 197
9 Wq = torch.randn(d_q, d)
10 Wk = torch.randn(d_q, d)
11 X = torch.randn(d, num_tokens)
12 W = torch.mm(Wq.T, Wk)
13
14
15 # Normal attention mode
16 W1 = W
17 P = torch.mm(torch.mm(X.T, W1), X)
18 attn_map1 = P.softmax(dim=1)
19
20
21 # Malignant attention entropy collapse mode
22 u,s,v = torch.svd(W)
23 s[0:3] = torch.tensor([3., 2., 1.])*s[0:3]
24 s[3:] = 0.0
25 W2 = torch.mm(torch.mm(u, torch.diag(s)), v.T)
26 P = torch.mm(torch.mm(X.T, W2), X)
27 attn_map2 = P.softmax(dim=1)
28
29
30 # Benign attention entropy collapse mode
31 u,s,v = torch.svd(W)
32 W3 = torch.mm(torch.mm(u, torch.diag(s)), u.T)
33 P = torch.mm(torch.mm(X.T, W3), X)
34 attn_map3 = P.softmax(dim=1)
35
36
37 # Plot figures
38 fig, axs = plt.subplots(1, 3, figsize=(15, 5))
39 axs[0].imshow(attn_map1.detach().numpy())
40 axs[1].imshow(attn_map2.detach().numpy())
41 axs[2].imshow(attn_map3.detach().numpy())
42 plt.show()

```

H ATTENTION MAP VISUALIZATION OF GPT

Figure 7 visualizes the dynamic process of attention map as the number of training steps increases for a successful and unsuccessful GPT-Small model. It should be noted that the GPT model uses a lower triangular attention mask.

(a) Block 11 (successful). (b) Block 11 (unsuccessful).

FIGURE 7: Visualization of the dynamic process of attention map as the number of training steps increases for a successful and unsuccessful GPT-Small model. Attention map gradually becomes sparse and low-rank along with the training process in a failure case. *Please click the images to play the flash. Best viewed with Acrobat Reader.*

In Figure 7, the attention values in a successful case distribute to different position, but the attention values in a unsuccessful case will only concentrate into several directions.

I MORE TRAINING DYNAMICS OF ViT AND GPT

Figure 8 visualizes a successful ViT training process. Compared with Figure 1, we find several significant differences as follows.

- In a successful ViT training process, the value of $\sigma_1(\mathbf{W}_q^\top \mathbf{W}_k)$ increases to 16,000, then starts to oscillate smoothly. But for an unsuccessful training, the value suddenly increases to a very large value, around 300,000, it triggers the model crash,
- The γ_1 and β_1 in a successful ViT training process are very smooth, but they change a lot in an unsuccessful case,
- The fast increase of $\sigma_1(\mathbf{W}_q^\top \mathbf{W}_k)$ is accompanied by a fast increase of \mathbf{W}_q and \mathbf{W}_k .

We can observe similar phenomenon in Figure 9 and Figure 15. In a successful GPT training process, the value of $\sigma_1(\mathbf{W}_q^\top \mathbf{W}_k)$ increases to 60, then starts to oscillate smoothly. But for an unsuccessful GPT training, the value increases to 20,000. The difference between the scale of value between GPT and ViT may be due to the density and sparsity of the supervision signal. In GPT, each token will contribute a gradient, but in ViT, only one class label in an image provides a supervision information.

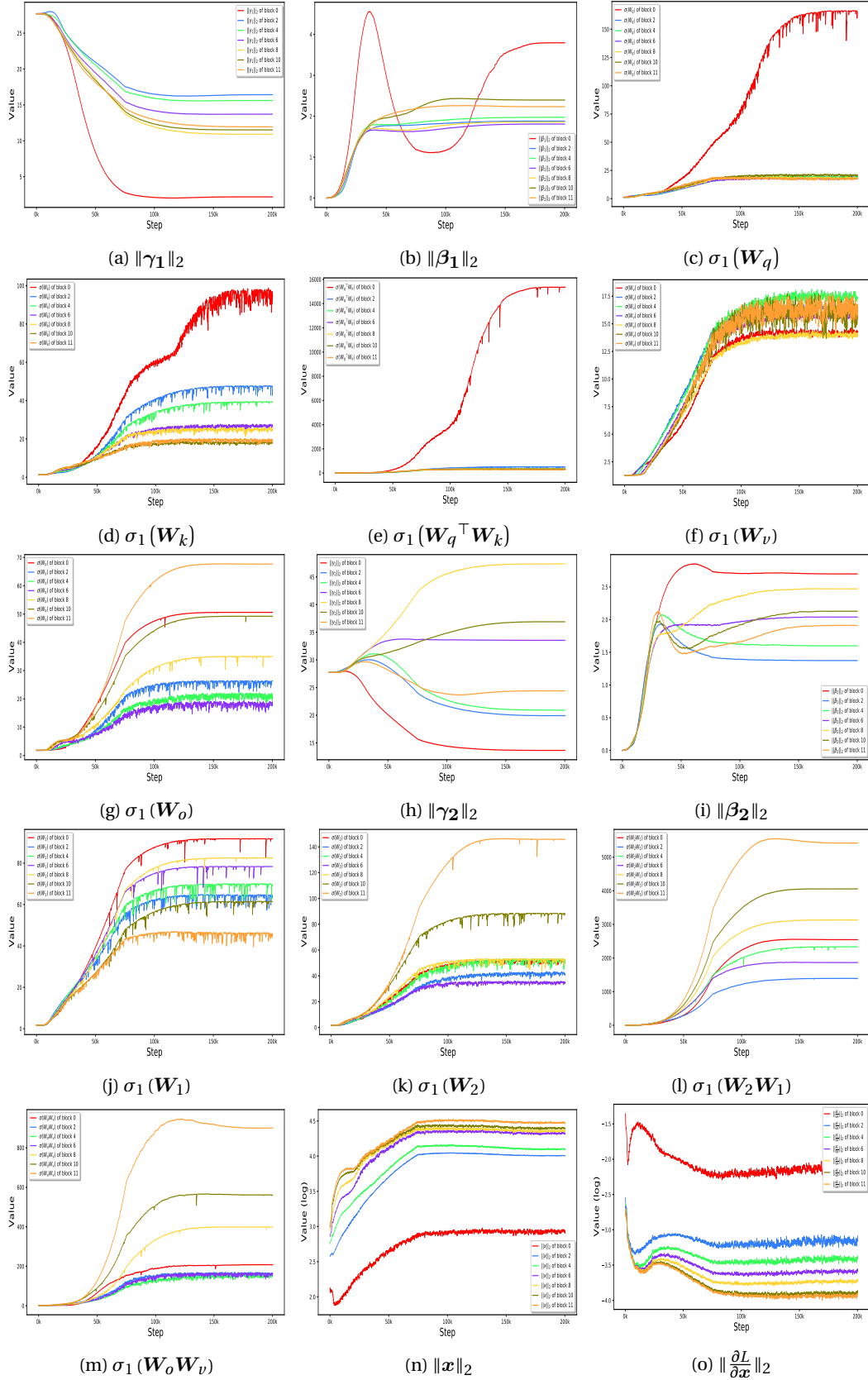


FIGURE 8: Training dynamics of a successful ViT training.

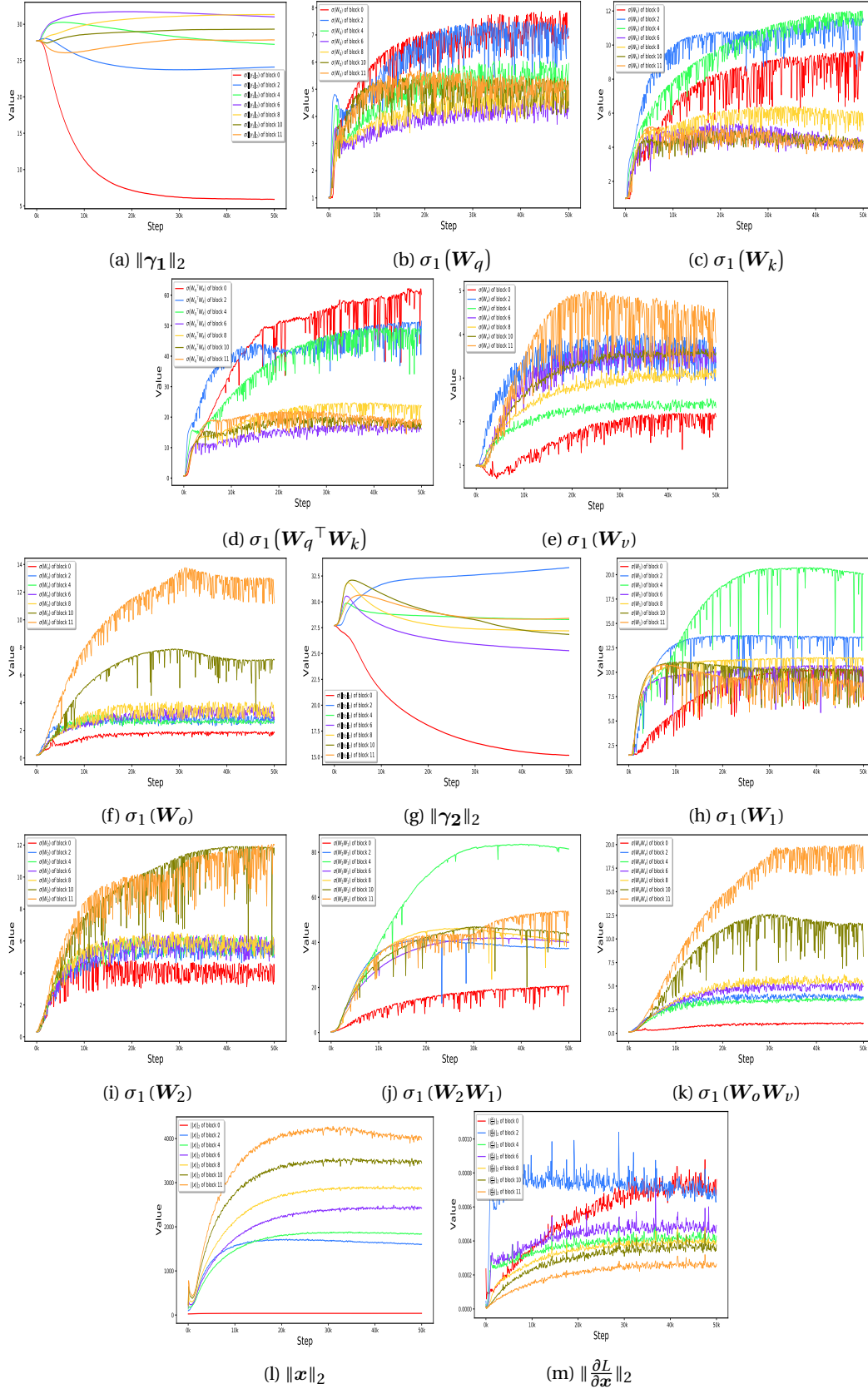


FIGURE 9: Training dynamics of a successful GPT training.

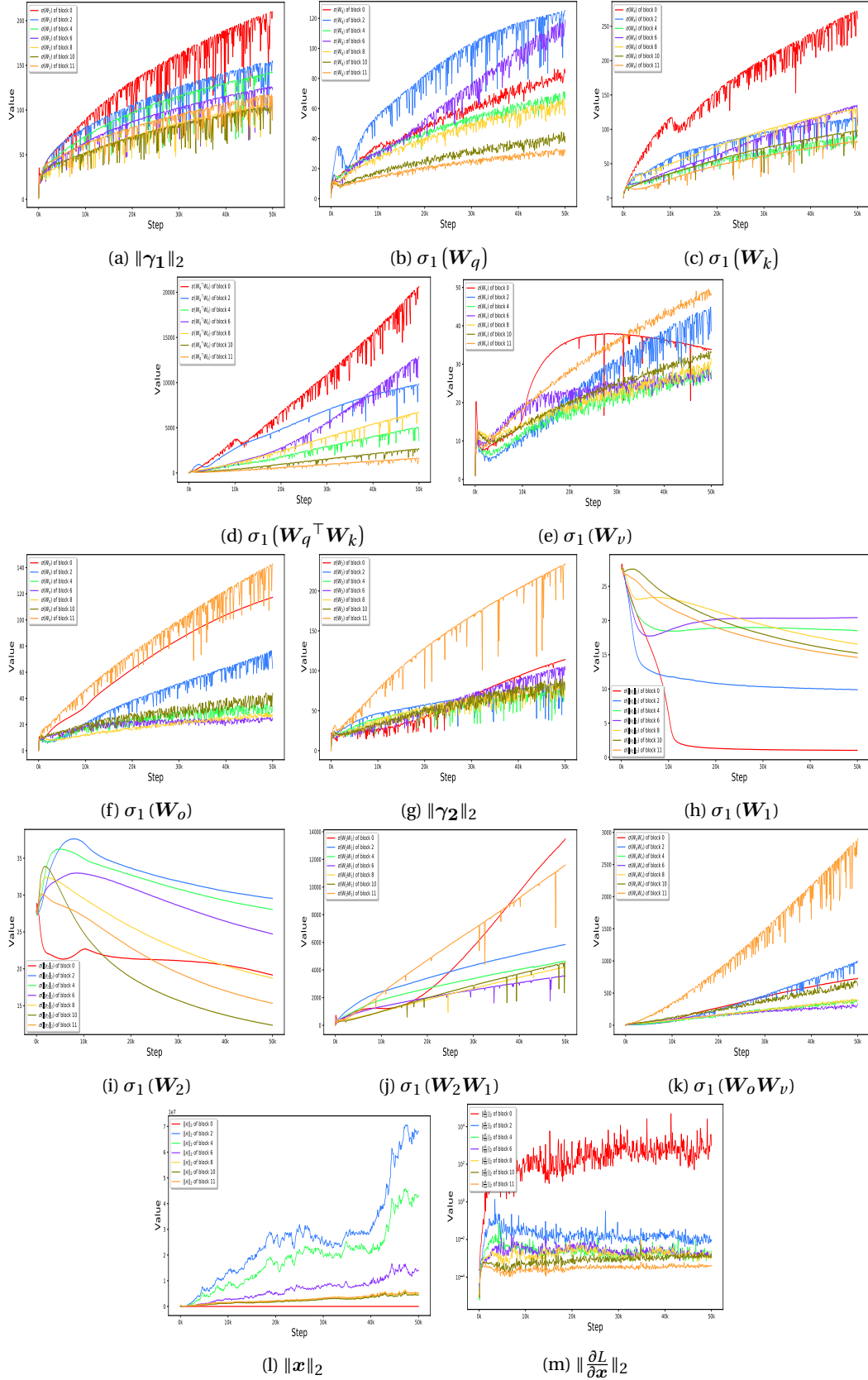


FIGURE 10: Training dynamics of an unsuccessful GPT training.

J EXPERIMENT OF 1B ViT

To further evaluate the effectiveness of our method at a larger scale, we assessed ViT-g with 1B parameters. The ViT-g 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 ViT-g with AdamW² and ViT-g with AdamW, where ViT-g with AdamW was evaluated under two settings: with and without learning rate warmup. Our AdamW² does not use warmup. The comparison results are presented in Figure 11 and Figure 12.

Figure 11 shows that ViT-g with AdamW crashes after only a few training steps when running without warmup. While the use of warmup enables ViT-g to complete training, but the loss spikes one time. Our ViT-g with AdamW² not only achieves stable training without warmup but also demonstrates better performance.

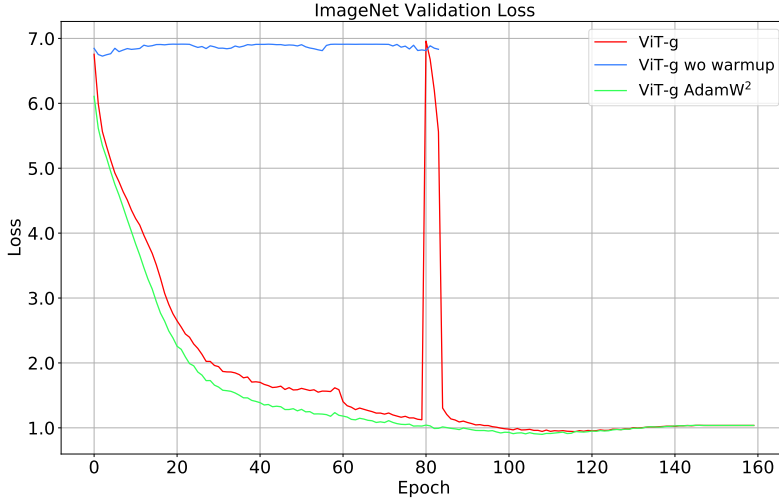


FIGURE 11: Comparison of loss curve of AdamW² and AdamW on ViT-g model.

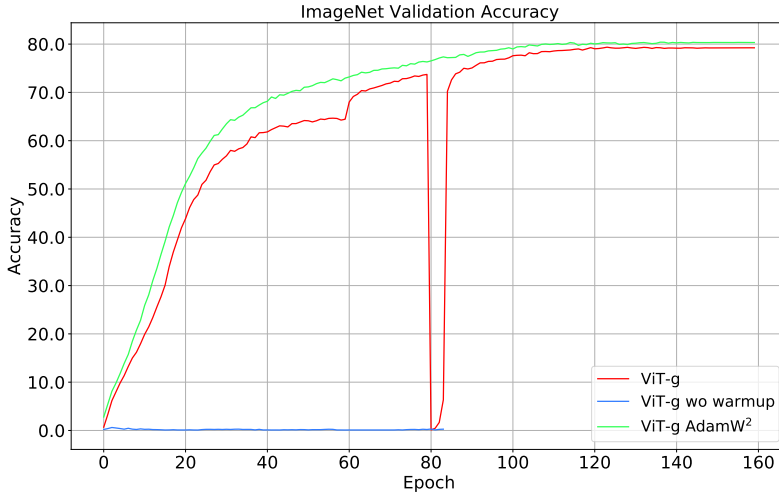


FIGURE 12: Comparison of accuracy of AdamW² and AdamW on ViT-g model.

K EXPERIMENT OF 774M NANO GPT

We also evaluated the effectiveness of our method on a larger-scale language model, termed as nanoGPT-large. The model architecture consists of 36 layers with a hidden dimension of 1280 and 20 attention heads. The total parameter count is 774M. Our experimental setup strictly follows the nanoGPT configuration, including all learning rate settings. It is important to note that training nanoGPT-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 13. We can see from Figure 13, nanoGPT-large achieves a stable training without warmup and obtains a similar validation loss with its counterpart, GPT2-large. This further verifies our understanding to the model crash of Transformer.

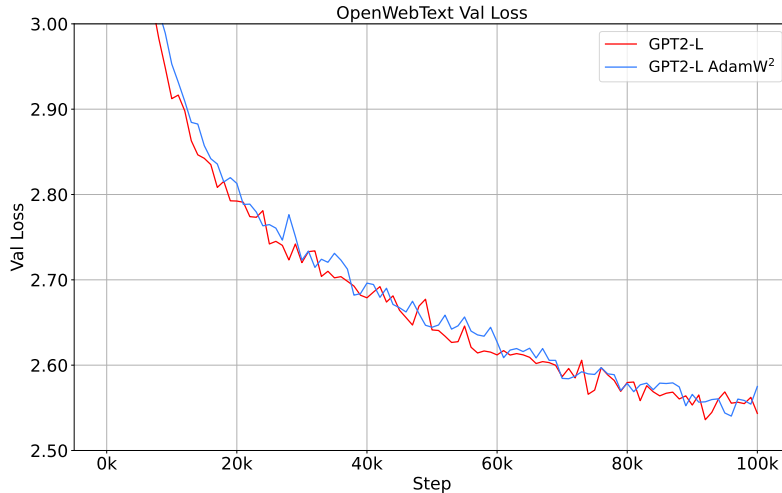


FIGURE 13: Comparison of validation loss of AdamW² and AdamW on nanoGPT-large model.

L EXPERIMENT OF FLATTEN-SWIN

Besides ViT, GPT, and Swin-Transformer, we further validated our approach on Flatten-Transformer Han et al. (2023). We used Flatten-Swin, and we compared the performance of our method and the baseline method training 150 and 300 epochs. Our method could stably train and demonstrate performance comparable to the baseline. This further verified the correctness of our understanding of neural network stability.

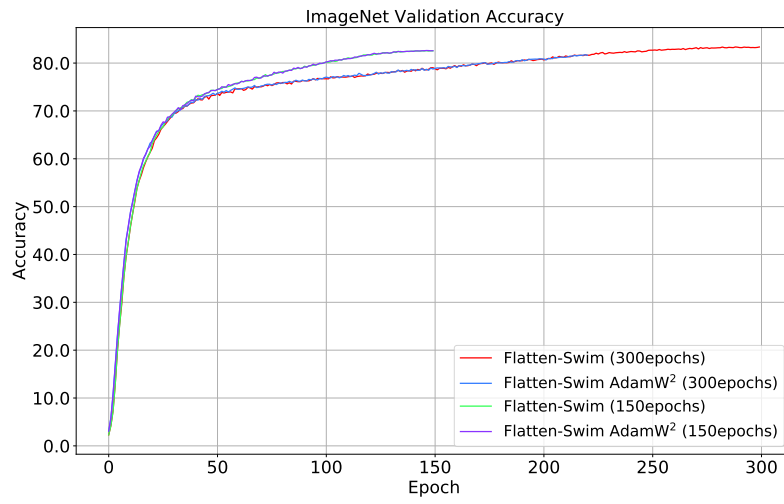


FIGURE 14: Evaluation of Flatten-Swin.

M ACTUAL LEARNING RATE CURVE ALONG WITH TRAINING STEPS

We recorded the actual learning rate throughout the training steps, we sample one point every 50 steps. Our initial setting of the learning rate is a cosine learning rate scheduler without warmup. If $\alpha_t \frac{\sigma_1(\nabla \mathbf{W}_t)}{\sigma_1(\nabla \mathbf{W}_{t-1})} > \tau$, then α_t will be truncated to $\tau \frac{\sigma_1(\nabla \mathbf{W}_{t-1})}{\sigma_1(\nabla \mathbf{W}_t)}$. From the figure, we observe that γ_1 and γ_2 in RMSNorm only exceed the preset τ during the initial training phase and rarely exceed it afterwards. For other curves, they somewhat look like a curve with learning rate warmup, but we can see that different blocks have different learning rates.

We also observe that shallower layers are more likely to violate the preset τ value. It means the shallower layers are more likely to lead to a greater update of weight matrix and typically require a smaller learning rate. Additionally, we notice that for the weight matrix \mathbf{W}_2 , it is more prone to exceeding the preset τ value compared to the weight matrix \mathbf{W}_1 .

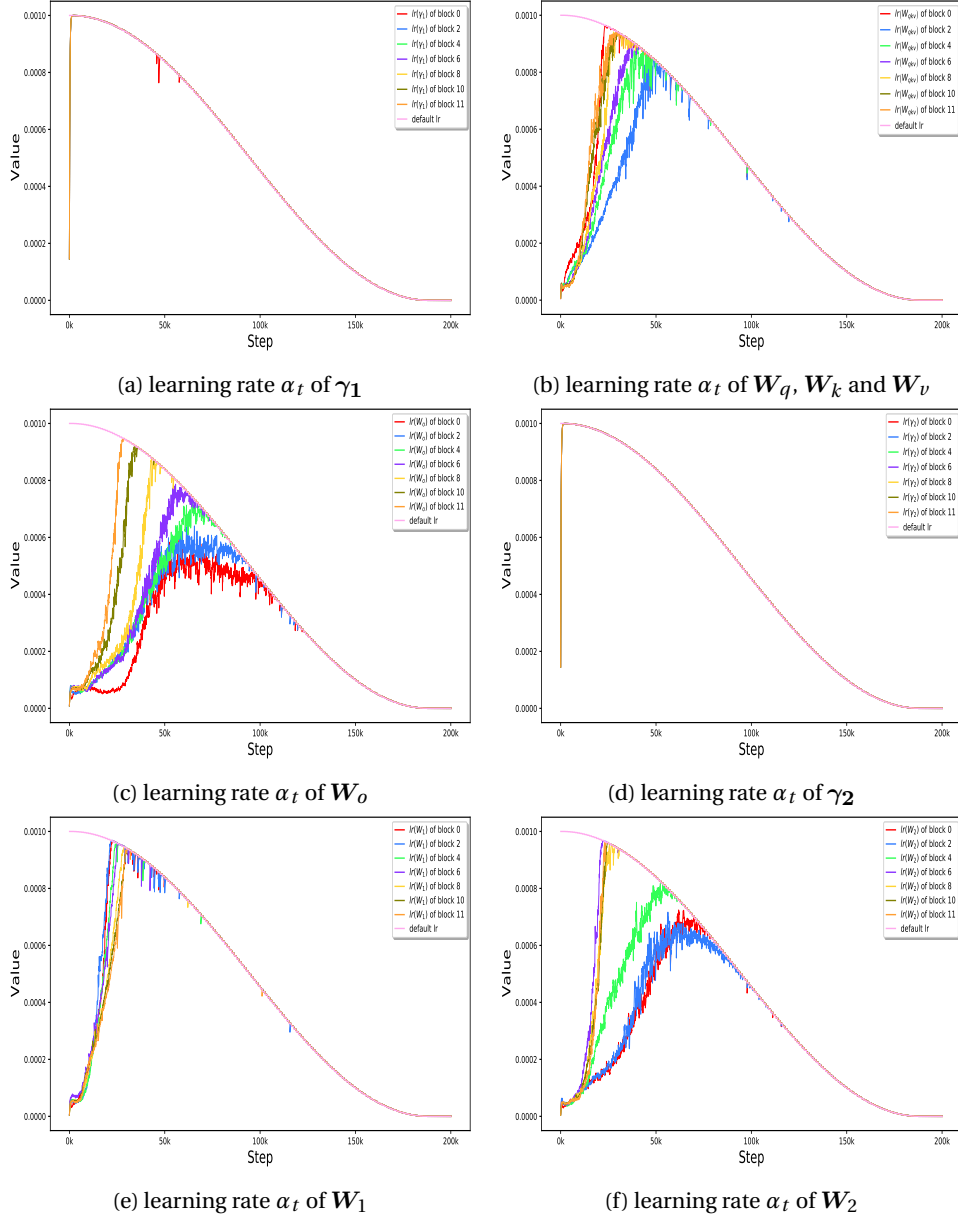


FIGURE 15: Actual Learning Rate Curve along with Training steps.

N TRAINING CONFIGURATIONS

Training Configurations. We list the training configurations of ViT, GPT, Swin-Transformer and Flatten-Swin in Table 2. For ViT, GPT, Swin-Transformer and Flatten-Swin, we do not use learning rate warmup. For GPT, we follow the experimental configurations of nanoGPT (Karpathy, 2022), all parameters are same as GPT2 (Radford et al., 2019). For ViT, we use Timm (Wightman, 2019). For Swin-Transformer, we use the original code provided by Liu et al. (2021). For Flatten-Swin, we use the original code provided by (Han et al., 2023).

TABLE 2: Training configurations for ViT, GPT and Swin-Transformer.

(a) Training configurations for ViT.		(b) Training configurations for GPT.	
training config	ViT-B/L/g (224 ²)	training config	GPT-S/L
optimizer	AdamW ²	optimizer	AdamW ²
τ (In default)	0.004 or 0.003	τ	0.01
warmup epochs	0	warmup epochs	0
weight init	Truncated Xavier	weight init	Xavier
base learning rate	1e-3	baseline learning rate	0.0006 or 0.00025
weight decay	0.05/0.1	weight decay	0.1
optimizer momentum	$\beta_1, \beta_2 = 0.9, 0.99$	optimizer momentum	$\beta_1, \beta_2 = 0.9, 0.95$
batch size	1024	tokens seen each update	500,000
training epochs	150	max iters	600K or 100K
learning rate schedule	cosine decay	batch size	480
randaugment	(9, 0.5)	sequence length	1024
mixup	0.8	dropout	0.0
cutmix	1.0	bfloat16	True
random erasing	0	gradient clipping	1.0
label smoothing	0.1		
stochastic depth	0.1/0.5		
gradient clip	None		
exp. mov. avg. (EMA)	no		

(c) Training configurations for Swin-Transformer.		(d) Training configurations for Flatten-Swin.	
training config	Swin S/B (224 ²)	training config	Flatten-Swin S (224 ²)
optimizer	AdamW ²	optimizer	AdamW ²
τ (In default)	0.004	τ (In default)	0.004
warmup epochs	0	warmup epochs	0
training epochs	300	training epochs	150 or 300
others	same as Liu et al. (2021)	others	same as Han et al. (2023)

O NON-SYMMETRIC POSITIVE QUASI-DEFINITE SQUARE MATRIX

When we mention a non-symmetric positive quasi-definite square matrix, we mean it has the following three properties,

1. $\mathbf{W}_q^\top \mathbf{W}_k$ is not symmetric because generally, $\mathbf{W}_q^\top \mathbf{W}_k \neq \mathbf{W}_k^\top \mathbf{W}_q$,
2. $\mathbf{W}_q^\top \mathbf{W}_k$ is a square matrix and most of its eigenvalues are larger than 0, and only very few are less than 0.0. So we call it positive quasi-definite matrix.
3. if we assume $\mathbf{W} = \mathbf{W}_q^\top \mathbf{W}_k$ is positive definite matrix, if for each element in \mathbf{x} is sampled from a standard Gaussian distribution, we can prove

$$\mathbb{E}[\mathbf{x}_i^\top \mathbf{W} \mathbf{x}_i] \gg \mathbb{E}[\mathbf{x}_i^\top \mathbf{W} \mathbf{x}_j]$$

when $i \neq j$, see Appendix C for the proof.

P DISCUSSION ABOUT RANK COLLAPSE, ENTROPY COLLAPSE AND SPARSE YET LOW-RANK ENTROPY MATRIX

Before we start our discussion, let us see three matrices,

$$\mathbf{A} = \begin{pmatrix} 1 & 1 & 1 & 1 & 1 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \vdots \end{pmatrix}, \mathbf{B} = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 \end{pmatrix}, \mathbf{C} = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 \end{pmatrix}$$

We can see that \mathbf{A} is low-rank, \mathbf{B} is sparse but not low-rank, \mathbf{C} is sparse and low-rank.

In previous papers (Dong et al., 2021; Zhai et al., 2023), researchers have analyzed the problem of model crash via *rank collapse of activations and entropy collapse of attention map*. Dong et al. (Dong et al., 2021) attributes the model crash into *rank collapse of the activations*, but Zhai et al. (2023) think it is the *entropy collapse of the attention map* leading to the model crash.

However, based on our analysis, we can find a counterexamples for entropy collapse, and meanwhile the rank collapse of the activation cannot fully describe the inner reason of the model crash (the weight matrix instead of activations). When the state of \mathbf{C} usually happens, the model crashed,

- Rank collapse of the activations cannot reveal the underlying cause that exists in the weight matrix. Weight matrix is the inner key ingredient of the model instead of activations.
- \mathbf{B} is a counterexample of entropy collapse. we observe that in some successful cases, state \mathbf{B} occurs. According to the definition of entropy collapse, state \mathbf{B} should lead to model crash; however, our experiments show that the model remains stable in this state.
- Sparse yet low-rank attention matrix is the state of the attention map when a model crashes. We believe rank collapse of activations and entropy collapse of attention map are not enough to describe the state of the model crash precisely. According to our analysis, the Spectral Energy Concentration (SEC) of the $\mathbf{W}_q^\top \mathbf{W}_k$ is the inner reason the model crash, and the sparse yet low-rank attention matrix is the phenomena observed on the attention matrix.

In summary, our paper, via a rigid theoretical analysis, our paper reveals the Spectral Energy Concentration (SEC) of the $\mathbf{W}_q^\top \mathbf{W}_k$ is the inner reason the model crash, and the sparse yet low-rank attention matrix is the phenomena that is observed on the attention matrix.

Q ILLUSTRATION OF FIGURE 5

Illustration of Arrow 1.

According to the property of Kronecker Product, we have

$$\text{rank}(\mathbf{X} \otimes \mathbf{X}) = \text{rank}(\mathbf{X}) \cdot \text{rank}(\mathbf{X}).$$

Since \mathbf{X} is low-rank, then $\mathbf{X} \otimes \mathbf{X}$ is also low-rank. In the following, we will also prove the singular values of $\mathbf{X} \otimes \mathbf{X}$ will also strengthen the concentration of spectral energy into some directions with large singular values.

Illustration of Arrow 2.

According to the computation of Jacobian matrix, we have

$$\frac{\partial \text{vec}(\mathbf{P})}{\partial \text{vec}(\mathbf{W}_q^\top \mathbf{W}_k)} = \mathbf{X}^\top \otimes \mathbf{X}^\top$$

where $\mathbf{P} = \mathbf{X}^\top \mathbf{W}_q^\top \mathbf{W}_k \mathbf{X}^\top$.

Illustration of Arrow 3.

Let $\mathbf{X} \in \mathbb{R}^{m \times n}$ be a matrix with rank $r \leq \min(m, n)$. Denote the singular values of \mathbf{X} as $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0$ and $\sigma_{r+1} = \dots = \sigma_{\min(m, n)} = 0$.

Definition. (Singular Values of Kronecker Product)** For a matrix \mathbf{X} , define $\Lambda(\mathbf{X})$ as the set of all possible products of its singular values, i.e., $\Lambda(\mathbf{X}) = \{\sigma_i \sigma_j : 1 \leq i, j \leq r\}$.

Theorem (Singular Values of Kronecker Product) For a low-rank matrix $\mathbf{X} \in \mathbb{R}^{m \times n}$ with rank r , the singular values of $\mathbf{X} \otimes \mathbf{X}$ are precisely the elements in $\Lambda(\mathbf{X})$. More formally, let $\{\mu_k\}$ be the set of singular values of $\mathbf{X} \otimes \mathbf{X}$. Then $\{\mu_k\} = \{\sigma_i \sigma_j \text{ where } 1 \leq i, j \leq r\}$.

Proof. Consider the singular value decomposition (SVD) of $\mathbf{X} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^\top$, where $\mathbf{U} \in \mathbb{R}^{m \times m}$ is an orthogonal matrix, $\mathbf{V} \in \mathbb{R}^{n \times n}$ is an orthogonal matrix, and $\mathbf{\Sigma} = \text{diag}(\sigma_1, \dots, \sigma_r, 0, \dots, 0)$. Then, the Kronecker product $\mathbf{X} \otimes \mathbf{X}$ can be expanded as: $\mathbf{X} \otimes \mathbf{X} = (\mathbf{U} \mathbf{\Sigma} \mathbf{V}^\top) \otimes (\mathbf{U} \mathbf{\Sigma} \mathbf{V}^\top) = (\mathbf{U} \otimes \mathbf{U})(\mathbf{\Sigma} \otimes \mathbf{\Sigma})(\mathbf{V}^\top \otimes \mathbf{V}^\top)$. The singular values of $\mathbf{U} \otimes \mathbf{U}$ and $\mathbf{V}^\top \otimes \mathbf{V}^\top$ are all 1, as they are composed of orthogonal matrices.

Therefore, the singular values of $\mathbf{X} \otimes \mathbf{X}$ are precisely the elements of $\mathbf{\Sigma} \otimes \mathbf{\Sigma}$, which are exactly $\{\sigma_i \sigma_j : 1 \leq i, j \leq r\}$.

In summary, the Singular Values of Kronecker Product $\mathbf{X} \otimes \mathbf{X}$ will also strengthen the concentration of spectral energy into some directions with large singular values. In this way, the singular values of $\frac{\partial \text{vec}(\mathbf{P})}{\partial \text{vec}(\mathbf{W}_q^\top \mathbf{W}_k)}$ will also concentrate to a few directions.

When we update $\mathbf{W}_q^\top \mathbf{W}_k$ according to the following equation:

$$\text{vec}(\mathbf{W}_q^\top \mathbf{W}_k)_{\text{new}} = \text{vec}(\mathbf{W}_q^\top \mathbf{W}_k)_{\text{old}} - \alpha \frac{\partial \mathcal{L}}{\partial \text{vec}(\mathbf{P})} \frac{\partial \text{vec}(\mathbf{P})}{\partial \text{vec}(\mathbf{W}_q^\top \mathbf{W}_k)}$$

where \mathcal{L} is the loss function, and α is the step size. Since the singular values of $\frac{\partial \text{vec}(\mathbf{P})}{\partial \text{vec}(\mathbf{W}_q^\top \mathbf{W}_k)}$ will also concentrate to a few directions, the update will lead to the singular values of $\mathbf{W}_q^\top \mathbf{W}_k$ tends to concentrate.

Illustration of Arrow 4.

Our Theorem 1 in the paper is to demonstrate that the spectral energy concentration of $\mathbf{W}_q^\top \mathbf{W}_k$ and the associated dominant large singular values will lead to \mathbf{A} to be a sparse yet low-rank matrix.

Illustration of Arrow 5.

Since $\mathbf{X}^{l+1} = \mathbf{V} \mathbf{A}$, where \mathbf{V} is the value matrix after projection in attention module and \mathbf{A} is the attention matrix, according to linear algebra, we have,

$$\text{rank}(\mathbf{X}^{l+1}) < \min\{\text{rank}(\mathbf{V}), \text{rank}(\mathbf{A})\}.$$

Thus, we have that \mathbf{X}^{l+1} is also a low-rank matrix.