

# ONLY LARGE WEIGHTS (AND NOT SKIP CONNECTIONS) CAN PREVENT THE PERILS OF RANK COLLAPSE

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Attention mechanisms lie at the heart of modern large language models (LLMs). Straightforward algorithms for forward and backward (gradient) computation take quadratic time, and a line of work initiated by [Alman and Song NeurIPS 2023] and [Alman and Song NeurIPS 2024] has shown that quadratic time is necessary unless the model weights are small, in which case almost linear time algorithms are possible. In this paper, we show that large weights are necessary to avoid a strong preclusion to representational strength we call layer collapse, which means that the entire network can be approximated well by a network with only a single layer. This means that transformers with small weights are shockingly weak, and that the quadratic running time of attention is unavoidable for expressive transformers.

The notion of layer collapse that we introduce is a variant on the notion of rank collapse from the work of [Dong, Cordonnier, and Loukas ICML 2021]. They showed that in Self Attention Networks with small weights and with skip connections, rank collapse must occur. This is typically interpreted as justifying the necessity of skip connections in expressive networks. However, our result shows that even with skip connections, if the weights are small, then layer collapse still occurs. Thus, only large weights, and not skip connections, can prevent these representational weaknesses.

## 1 INTRODUCTION

The rapid progress of large language models, text-to-image and text-to-video models like Transformer (Vaswani et al., 2017), BERT (Devlin et al., 2018), GPT-4 (OpenAI, 2023), Llama 3 (Llama Team, 2024), and Gemini 2.0 (Google, 2025), has enabled powerful language modelling abilities. These models take advantage of large-scale pretraining on massive textual data, which equips them with strong abilities to interpret the complex patterns of natural language. These LLMs have a broad range of applications, influencing domains such as human-computer interaction, multilingual translation, language comprehension, text generation, and rapid prototyping of software.

The major architecture behind the success of all these language models is the attention mechanism. Specifically, attention computes pairwise similarities by calculating inner products between vectorized representations of words, with input sequences represented as vectors. Formally, softmax attention can be formulated as follows:

**Definition 1.1** (Self-Attention with Softmax Units). *Let  $A \in \mathbb{R}^{n \times d}$  and weights  $Q, K, V \in \mathbb{R}^{d \times d}$ . Let  $g$  represent the entry-wise exponentiation function, i.e., for  $z \in \mathbb{R}$  we have  $g(z) = \exp(z)$ , and for a matrix  $W$  we have  $g(W)_{i,j} = g(W_{i,j})$ . The attention computation can be defined as*

$$\text{SAtt}(X, Q, K, V) = \underbrace{D^{-1}}_{n \times n} \underbrace{g(XQK^\top X^\top)}_{n \times n} \underbrace{X}_{n \times d} \underbrace{V}_{d \times d}$$

where  $D := \text{diag}(g(XQK^\top X^\top)\mathbf{1}_n)$ , and where  $\mathbf{1}_n \in \mathbb{R}^n$  is a length- $n$  vector whose entries are all 1.

**Small Coefficients are Needed for Fast Algorithms** However, the straightforward algorithm for computing self-attention results in a quadratic  $O(n^2d)$  running time, where  $n$  is the length of the input token and  $d$  is the hidden dimension. Under popular complexity-theoretic assumptions, there is no better, subquadratic time algorithm to compute attention, even approximately (Alman & Song, 2023). Therefore, models based on attention may face difficulties when they handle long contexts.

In fact, a key observation of this line of work on the computational complexity of attention is that attention can be computed (or tightly approximated) faster if one restricts to small weights, i.e., an upper bound on how large the entries of  $Q, K, V$  can be in Definition 1.1 above. Indeed, a line of work (Alman & Song, 2023; 2024a;b; 2025) has shown that small weights are both necessary and sufficient for a faster algorithm: If the weights are large, then the aforementioned complexity-theoretic result shows that there is no subquadratic time algorithm. However, if the weights are small, then attention can be approximated to low error in *almost linear time*! Their algorithm is based on low-rank approximations of the  $n \times n$  attention matrix (the matrix  $g(XQK^\top X^\top)$  in Definition 1.1 above).

This type of observation is also frequently used in practice; many LLM implementations have enforced bounds on the weights, often using techniques like approximation or quantization, and then used this for substantial speedups. For some examples, see (Zafir et al., 2019; Katharopoulos et al., 2020b; Frantar et al., 2022; Perez et al., 2023; Dettmers et al., 2023; Egashira et al., 2024; Liu et al., 2024b; Xu et al., 2024a; Lin et al., 2025; Chen et al., 2025b; Liu et al., 2025; Ouyang et al., 2025; Deng et al., 2025; Hu et al., 2025c; Fu et al., 2025; Hu et al., 2025b; Park et al., 2025; Zeng et al., 2025; Yu et al., 2025; Wei et al., 2025).

In this paper, we investigate the representational strength of transformers with small weights. Our main result will show a limitation, that without large weights, a transformer cannot take advantage of more than a single layer. In other words, we will show that in order to take advantage of the full expressive power of the transformer model, large weights are necessary.

**Rank Collapse and Skip Connections** We will crucially build on the approach of Dong, Cordonnier, and Loukas (Dong et al., 2021), who studied the representational strength of different variants on the transformer architecture through the lens of a notion called *rank collapse*. We say that a model experiences rank collapse if, on any input, the output must always be close to a rank 1 matrix. (See Definition 3.4 below for the precise meaning.) Beyond being unable to represent complex concepts, models with rank collapse also have numerous other issues in both training and evaluation (Noci et al., 2022; Roth & Liebig, 2024; Naderi et al., 2024; Nguyen et al., 2024; Heo & Choi, 2024; Yuan & Xu, 2024; Barbero et al., 2025; Bonino et al., 2025).

The work of (Dong et al., 2021) highlights *skip connections* (or residual connections) in a transformer network as crucial for avoiding rank collapse. They show that in a Self-Attention Network without skip connections, rank collapse occurs with a doubly exponential rate of convergence. More precisely, if  $\beta$  is a bound on the  $\ell_1$  norm of the weight matrices of the network, and the network has  $L$  layers, then they show the distance to a rank-1 matrix shrinks as

$$O(\beta)^{\frac{3^L-1}{2}}. \quad (1)$$

Meanwhile, they observe that networks with skip connections may experience no rank collapse at all. For instance, it is not hard to simulate the *identity* function as a Self-Attention Network with skip connections (simply set all value weights to 0, so that only the skip connections are output). In this case, any input which is far from rank-1 will result in an output which is also far from rank-1. They study other mechanisms in transformer networks as well, including multi-layer perceptrons and layer normalization, but find that only skip connections prevents the rank collapse of Equation (1). This result is frequently cited in the literature as evidence of the importance of skip connections (Ma et al., 2021; Noci et al., 2022; Sander et al., 2022; Guo et al., 2023; Li et al., 2023; Kim et al., 2023; Geshkovski et al., 2023; Kim et al., 2024; Ji et al., 2025).

**The Importance of Large Weights and Layer Collapse** We begin with a simple observation: in order for Equation (1) to be shrinking as  $L$  grows, it is necessary that  $\beta$  is small, i.e., that the weights of the network are small. In other words, the result of (Dong et al., 2021) really says that:

To avoid rank collapse, one needs either skip connections *or* large weights.

In this paper, we prove that Self Attention Networks with skip connections, but with small weights, must suffer from a phenomenon similar to rank collapse which we call *layer collapse*. We say that an  $L$ -layer Self Attention Network  $S$  has layer collapse if there is a nearly equivalent Self Attention Network  $S'$  which only has a single layer. In other words, although  $S'$  only has one layer, it is still as expressive as  $S$ , since on any input  $X$ , the outputs  $S(X)$  and  $S'(X)$  differ in each entry by at most a small error parameter.

When combined with (Dong et al., 2021), our result implies:

To avoid rank and layer collapse, one needs large weights (skip connections do not suffice).

This challenges the previous popular interpretation of (Dong et al., 2021), that skip connections were crucial for the representational strength of the model.

The connection between layer collapse and rank collapse may not be evident from the definitions, but it will become clear in our proofs below. At a high level, we will find that the attention mechanisms in lower layers of the Self Attention Network must exhibit rank collapse (regardless of skip connections), and can thus be removed from the network without substantially changing the output. We will show

**Theorem 1.2** (Main result, informal). *If  $S$  is a Self Attention Network whose weight matrices have  $\ell_\infty$  norm bounded by  $\eta$ , then there is a Self Attention Network  $S'$  with only one layer, such that on any input  $X$  with  $\|X\|_\infty \leq O(1)$ , we have  $\|S(X) - S'(X)\|_\infty \leq O(\eta)$ .*

In fact, the example from (Dong et al., 2021) of the identity network with skip connections heavily inspired our definition of layer collapse. That network indeed does not have rank collapse, so we could not hope to prove a version of Theorem 1.2 with rank collapse instead of layer collapse. On the other hand, it is essentially not making use of its attention mechanisms; they could be removed without changing the output of the network. Our key idea is to show that, more generally, the attention mechanisms with small weights can be removed from any Self Attention Network, with skip connections, without changing the output of the network by very much.

Our  $\eta$  in Theorem 1.2 is a bound on the  $\ell_\infty$  norm of the weight matrices (maximum magnitude of an entry), whereas the prior result in Eq. (1) above uses parameter  $\beta$ , which is a bound on the  $\ell_1$  norm (sum of magnitudes of all entries). Our  $\eta$  could thus be quite a bit smaller (by a factor of  $d^2$  for  $d \times d$  weight matrices), and there are thus networks without skip connections where (Dong et al., 2021) does not imply rank collapse (since  $\beta \gg 1$  is too big) but our Theorem 1.2 still implies layer collapse (since  $\eta \ll 1$  is smaller).

We also note that both our informal statement of Theorem 1.2 and our presentation of the main result of (Dong et al., 2021) in Eq. (1), are given assuming that the Self Attention Network has a constant number of heads and layers. The more complete statement in terms of the number of heads and layers is presented in Theorem D.1 in the appendix. Both results have modest assumptions on the relationships between  $\eta$  (or  $\beta$ ),  $\|X\|_\infty$ , and the numbers of heads and layers, and we emphasize that these assumptions are nearly identical in both results; see Remark D.3 for more details.

**Roamdap.** In Section 2, we present the related work. In Section 3, we introduce several basic notations and definitions. In Section 4, we study perturbation properties of several functions, such as softmax. In Section 5, we provide several major rank collapse results. In Section 6, we provide the conclusion of this paper.

## 2 RELATED WORK

**Low-rank Approximations** Low rank approximation is a fundamental topic in numerical linear algebra (Clarkson & Woodruff, 2013; Nelson & Nguyễn, 2013; Song et al., 2023b;a). Many problems require either computationally or analytically finding a low-rank approximation under different settings such as linear and kernel SVMs (Gu et al., 2025), tensor regression (Song et al., 2021b; Reddy et al., 2022; Diao et al., 2018; 2019), low rank approximation with Frobenious norm (Clarkson & Woodruff, 2013; Nelson & Nguyễn, 2013), weighted low rank approximation (Razenshteyn et al., 2016; Gu et al., 2024; Li et al., 2025a; Song et al., 2025), general norm column subset selection (Song et al., 2019a), entrywise  $\ell_1$  norm low rank approximation (Song et al., 2017; 2019b),

162 tensor low rank approximation (Song et al., 2019c), tensor power method (Deng et al., 2023b), and  
163 matrix CUR decomposition (Boutsidis & Woodruff, 2014; Song et al., 2017; 2019c). Rank collapse  
164 and other techniques we use here build on this line of work.

165  
166 **Algorithmic Result for Attention Computations** The quadratic time complexity of attention  
167 mechanisms (Vaswani et al., 2017) has posed significant computational challenges for long se-  
168 quences. In response to this problem, a wide range of works have been proposed to reduce com-  
169 putational cost and enhance the scalability of attention mechanisms, including sparsification (Child  
170 et al., 2019; Zaheer et al., 2020; Beltagy et al., 2020; Hubara et al., 2021; Shi et al., 2023a; Kurtic  
171 et al., 2023; Frantar & Alistarh, 2023; Li et al., 2024b; Liang et al., 2024a; Han et al., 2024), kernel-  
172 based approaches (Liu & Zenke, 2020; Charikar et al., 2020; Zandieh et al., 2023; Deng et al., 2023a;  
173 Liang et al., 2024b), and low-rank methods (Li et al., 2016; Razenshteyn et al., 2016; Hu et al., 2022;  
174 2024b; Zeng & Lee, 2024). Additionally, another promising line of research is linear attention (Tsai  
175 et al., 2019; Katharopoulos et al., 2020a; Schlag et al., 2021; Deng et al., 2023c; Sun et al., 2023;  
176 Zhang et al., 2023b; Ahn et al., 2024; Li et al., 2024a; Shi et al., 2023c; Zhang et al., 2024), which  
177 significantly accelerates traditional softmax attention. Other relevant works have explored important  
178 aspects of attention mechanisms, covering topics such as circuit complexity (Chen et al., 2024a;c;  
179 Li et al., 2025b), model pruning (Frantar & Alistarh, 2023; Shen et al., 2024; Sun et al., 2024; Liang  
180 et al., 2025), privacy protection (Liang et al., 2024d; Gao et al., 2024), regression (Gao et al., 2023b),  
181 half-space reporting (HSR) (Jiang et al., 2021; Chen et al., 2024b), and quantum computation (Gao  
182 et al., 2023c; Zhao et al., 2024).

183 **Polynomial Kernels for Attention Acceleration** With the assumption that model weights are  
184 small, polynomial kernels (Alman & Song, 2023; 2024b) are powerful tools for approximating at-  
185 tention computation in almost linear time complexity, providing promising acceleration for both  
186 training and inference of a single attention layer. This approach can be further extended to a wide  
187 range of applications. For instance, polynomial kernels can provide insights into novel attention  
188 mechanisms and model designs, such as modern Hopfield models (Hu et al., 2024a), Diffusion  
189 Transformers (DiTs) (Shen et al., 2025; Hu et al., 2024d), multi-layer Transformers (Liang et al.,  
190 2024c), and tensor attention mechanisms (Liang et al., 2024e; Alman & Song, 2024a). These poly-  
191 nomial kernel methods also contribute to efficient and model-utility-preserving fine-tuning of foun-  
192 dation models, such as model adapters (Hu et al., 2022; Zhang et al., 2023a; Shi et al., 2023b; Gao  
193 et al., 2023a), multi-task fine-tuning (Gao et al., 2021; Oswald et al., 2023; Xu et al., 2024b), black-  
194 box model tuning (Sun et al., 2022), and instruction tuning (Li & Liang, 2021; Chung et al., 2022;  
195 Mishra et al., 2022). Other promising applications include privacy protection in attention computa-  
196 tion (Liang et al., 2024d), CoT reasoning (Khattab et al., 2022; Wei et al., 2022; Yao et al., 2023;  
197 Zheng et al., 2024), and model calibration (Zhao et al., 2021; Zhou et al., 2023). Very recently,  
198 (Gupta et al., 2025) further extends the work of (Alman & Song, 2023) to almost all the regimes of  
199 parameter  $d$  (see definition of  $d$  in Definition 1.1).

200 **Regression Models** The unprecedented energy consumption in training large-scale ML models  
201 has necessitated the development of scalable and efficient ML models (Venkataramani et al., 2015;  
202 Bender et al., 2021; McDonald et al., 2022). As a simple yet powerful approach to solving various  
203 machine learning problems (Bubeck, 2015; Brand et al., 2021; Song et al., 2024b; Subrahmanya &  
204 Shin, 2009), simple regression models have raised significant concerns in model acceleration, with  
205 recent advances from different perspectives, including sketching (Song & Yu, 2021; Reddy et al.,  
206 2022; Song et al., 2023a) and pre-conditioning (Yang et al., 2018; Kelner et al., 2022; Song et al.,  
207 2024a). Our work discusses low-rank approximations in attention mechanisms, while our general  
208 insight can be extended to other low-rank method applications, such as accelerated regression mod-  
209 els.

210 **Diffusion Models** Diffusion models and score-based generative models have achieved remark-  
211 able success in generating human-preference-aligned and high-quality visual content (Ho et al.,  
212 2020; Song et al., 2021a; Blattmann et al., 2023). These advances not only benefit vision tasks but  
213 also enhance the performance of other applications, such as language modeling (Lin et al., 2023;  
214 Sahoo et al., 2024), chemical design (Xu et al., 2023; Wen et al., 2024), and e-commerce (Yang  
215 et al., 2023; Wang et al., 2023; Liu et al., 2024a). Relevant works have discussed the theoretical  
guarantee that diffusion models can be approximated efficiently (Hu et al., 2024d; 2025a; 2024c;

Gong et al., 2025). Empirical approaches to accelerate diffusion models have addressed various aspects, such as shortcuts (Frans et al., 2024; Dao et al., 2024; Chen et al., 2025a), parameter pruning (Castells et al., 2024; Ma et al., 2024), and lazy computation (Nitzan et al., 2024; Shen et al., 2025). With these acceleration techniques, diffusion models can be trained on larger-scale data, overcoming inherent limitations such as counting (Hui et al., 2024; Cao et al., 2025; Guo et al., 2025a), text rendering (Chen et al., 2023; Tuo et al., 2024; Guo et al., 2025c), and adherence to physical constraints (Motamed et al., 2025; Guo et al., 2025b; Bansal et al., 2025). Most diffusion models leverage Transformer backbones for enhanced modelling capability. Our work accelerates attention mechanism computations, significantly benefiting a wide range of diffusion models.

**Graph ML Models** Relational data is prevalent in many real-world scenarios, where graph neural networks (GNNs) are the powerful solutions for mining effective patterns from such relations (Kipf & Welling, 2017; Hamilton et al., 2017; Wu et al., 2019). Recent scalability approaches have widely adopted low-rank approximations, such as sketching (Ding et al., 2022; Chamberlain et al., 2023) and vector quantization (Ding et al., 2021; Wang et al., 2025), which can take insights from this paper. These accelerations empower a wide range of applications, including misleading information mitigation (Xu et al., 2022; Chang et al., 2024), social network prediction (Fan et al., 2019; Zhang et al., 2022), and human action recognition (Peng et al., 2020; Li et al., 2021; Fu et al., 2021), while also inspiring advances in multiple aspects of graph learning, such as differential privacy (Lin et al., 2022; Mueller et al., 2022), robustness (Geisler et al., 2021; Dai et al., 2022; Zeng et al., 2022), and sensitive data removal (Chien et al., 2023; Zhang, 2024; Yi & Wei, 2025). A recent work (Zhang, 2024) proposes an efficient framework for empowering sensitive data impact removal from trained GNNs with partial retraining, leveraging model utility-aware data partitioning and contrastive sub-model aggregation.

### 3 PRELIMINARIES

In Section 3.1, we provide basic notation, definitions and facts. In Section 3.2 and Section 3.3, we define the Res function and balanced matrix notation which will appear prominently in our constructions. In Section 3.4, we provide the definition of a multi-layer multi-head Self Attention Network which we study here.

#### 3.1 BASIC NOTATION AND FACTS

For an arbitrary positive integer  $n$ , we use  $[n]$  to represent the set  $\{1, 2, \dots, n\}$ . We define  $\mathbf{1}_n$  as a length- $n$  vector where all entries are ones. For any  $x \in \mathbb{R}^n$ , we use  $\exp(x) \in \mathbb{R}^n$  to represent a length- $n$  vector whose  $i$ -th entry is  $\exp(x_i)$ . For any vector  $x \in \mathbb{R}^n$ , we use  $x^\top$  to denote its transpose. For a vector  $x$ , the vector  $\ell_2$  norm is denoted by  $\|x\|_2$ , i.e.,  $\|x\|_2 := (\sum_{i=1}^n x_i^2)^{1/2}$ . For a vector  $x$ , we use  $\|x\|_\infty$  to denote its  $\ell_\infty$  norm, i.e.,  $\|x\|_\infty := \max_{i=1}^n |x_i|$ . For a vector  $x$ , we use  $\|x\|_1$  to denote its entrywise  $\ell_1$  norm, i.e.,  $\|x\|_1 := \sum_{i=1}^n |x_i|$ . For a matrix, we use  $\|A\|_1$  to denote its  $\ell_1$  norm, i.e.,  $\|A\|_1 = \sum_{j,l} |A_{j,l}|$ . We use  $\|A\|_\infty$  to denote its  $\ell_\infty$  norm, i.e.,  $\|A\|_\infty := \max_{j,l} |A_{j,l}|$ . For a vector  $x \in \mathbb{R}^n$ , we use  $\text{diag}(x)$  to denote a diagonal matrix where  $i, i$ -th entry on diagonal is  $x_i$  for all  $i \in [n]$ .

**Definition 3.1.** For a vector  $x \in \mathbb{R}^n$ , we define  $\alpha(x) := \langle \exp(x), \mathbf{1}_n \rangle$ . We define  $\text{softmax}(x)$  as  $\text{softmax}(x) := \alpha(x)^{-1} \exp(x)$ . For a matrix  $A$ , we use the notation  $\text{softmax}(A)$  to denote that we apply  $\text{softmax}$  to each row of  $A$  individually. We define  $\text{softm}(X, W_Q, W_K)$  as follows

$$\text{softm}(X, W_Q, W_K) := \text{softmax}(XW_QW_K^\top X^\top) = \begin{bmatrix} \text{softmax}((XW_QW_K^\top X^\top)_{1,*}) \\ \text{softmax}((XW_QW_K^\top X^\top)_{2,*}) \\ \vdots \\ \text{softmax}((XW_QW_K^\top X^\top)_{n,*}) \end{bmatrix}$$

In many places, we will just write  $\text{softm}(X)$  for simplicity, and the weight matrices  $W_Q, W_K$  will be clear from context.

**Fact 3.2** (Shift-invariance property of softmax). For any vector  $x \in \mathbb{R}^n$  and for any fixed scalar  $a \in \mathbb{R}$ , we have  $\text{softmax}(x) = \text{softmax}(x + a\mathbf{1}_n)$ .

**Fact 3.3** (Norm inequality). For any matrices  $A \in \mathbb{R}^{n \times d}, B \in \mathbb{R}^{d \times m}$  we have (1)  $\|AB\|_1 \leq \|A\|_1 \cdot \|B\|_1$ , (2)  $\|AB\|_\infty \leq d\|A\|_\infty \cdot \|B\|_\infty$ , (3)  $\|AB\|_1 \leq m\|A\|_1 \cdot \|B\|_\infty$ .

### 3.2 DEFINITIONS OF Res

**Definition 3.4** (Res). Let  $Z \in \mathbb{R}^{n \times d}$  denote any matrix, we define function the Res :  $\mathbb{R}^{n \times d} \rightarrow \mathbb{R}^{n \times d}$  as  $\text{Res}(Z) := Z - \mathbf{1}_n y^\top$  where  $y := \arg \min_{y \in \mathbb{R}^d} \|Z - \mathbf{1}_n y^\top\|_\infty$ .

Res is the key definition behind the notion of rank collapse from prior work (Dong et al., 2021); we will use it here to study layer collapse as well, although we use the  $\infty$  norm here in contrast to prior work which uses a  $1, \infty$  norm.

### 3.3 $\theta$ -BALANCE

We also need a measure of how balanced a matrix is.

**Definition 3.5** ( $\theta$ -balance). Given a matrix  $E \in \mathbb{R}^{n \times n}$ , we define a corresponding matrix  $D \in \mathbb{R}^{n \times n}$  to be the diagonal matrix with  $D_{i,i} := \max_{j,l \in [n]} |E_{i,j} - E_{i,l}|$ . We say  $E$  is  $\theta$ -balanced, if  $\|D\|_\infty \leq \theta$ .

### 3.4 SELF-ATTENTION NETWORK

**Definition 3.6.** Let  $g$  denote the entry-wise exponentiation function, i.e., for  $z \in \mathbb{R}$  we have  $g(z) = \exp(z)$ , and for a matrix  $W$  we have  $g(W)_{i,j} = g(W_{i,j})$ . Given  $A \in \mathbb{R}^{n \times d}$  and weights  $Q, K, V \in \mathbb{R}^{d \times d}$ , the attention computation can be defined as

$$\text{SAtt}_H(X) := \sum_{h=1}^H \underbrace{D_h^{-1}}_{n \times n} \underbrace{g(XQ_h K_h^\top X^\top)}_{n \times n} \underbrace{X}_{n \times d} \underbrace{V_h}_{d \times d}$$

where  $D_h := \text{diag}(g(XQ_h K_h^\top X^\top) \mathbf{1}_n)$ , and where  $\mathbf{1}_n \in \mathbb{R}^n$  is a length- $n$  vector whose entries are all 1. When  $H = 1$  we simply write **SAtt** to denote the **SAtt<sub>H</sub>** function.

**Definition 3.7.** Let  $L, H$  denote fixed constants, where  $L$  represents the number of layers of the network, and  $H$  represents the number of heads per layer. Let **SAtt<sub>H</sub>** denote the multi-heads version of **SAtt** where  $H$  is the number of heads. For each  $\ell \in [L]$ ,  $X_\ell \in \mathbb{R}^{n \times d}$  denote the  $\ell$ -th layer input of self-attention network, then we have  $X_{\ell+1} = \text{SAtt}_H(X_\ell) + X_\ell$ .

## 4 PERTURBATION PROPERTY

We now move on to our main proof of layer collapse. We begin by showing that the relevant measure of matrices to not change much when their inputs are perturbed. We will ultimately show that layer collapse occurs because lower layers of the network can be seen as slightly perturbing their inputs. We study the Res function in Section 4.1, the  $\alpha$  function in Section 4.2.

### 4.1 PERTURBATION PROPERTY OF RES FUNCTION

**Lemma 4.1.** Let Res() be defined as Def. 3.4. If  $\|A - B\|_\infty \leq \epsilon$ , then  $\|\text{Res}(A) - \text{Res}(B)\|_\infty \leq 2\epsilon$ .

See Lemma E.1 in the Appendix for the proof of Lemma 4.1.

### 4.2 PERTURBATION PROPERTY OF EXP FUNCTION

**Lemma 4.2.** Let  $a, b \in \mathbb{R}^n$  such that  $\|b\|_\infty \leq \epsilon$ . Then, we can show: 1)  $|\exp(a_i + b_i) - \exp(a_i)| \leq (e^\epsilon - 1) \cdot \exp(a_i)$ . 2)  $|\exp(a_i + b_i) - \exp(a_i)| \leq (e^\epsilon - 1) \cdot \exp(a_i + b_i)$ . 3)  $|\alpha(a + b) - \alpha(a)| \leq (e^\epsilon - 1) \cdot \alpha(a)$ . 4)  $|\alpha(a + b) - \alpha(a)| \leq (e^\epsilon - 1) \cdot \alpha(a + b)$ .

*Proof.* It is easy to see that

$$\max\{|\exp(-b_i) - 1|, |\exp(b_i) - 1|\} \leq e^\epsilon - 1 \quad (2)$$

We can show

$$|\exp(a_i + b_i) - \exp(a_i)| = \exp(a_i) |\exp(b_i) - 1| \leq \exp(a_i) \cdot (e^\epsilon - 1) \quad (3)$$

where the first step follows from simple algebra, the second step follows from Eq. (2).

Thus, we have

$$\begin{aligned} |\alpha(a+b) - \alpha(a)| &= |\langle \exp(a+b), \mathbf{1}_n \rangle - \langle \exp(a), \mathbf{1}_n \rangle| \leq \sum_{i=1}^n |\exp(a_i + b_i) - \exp(a_i)| \\ &\leq \sum_{i=1}^n \exp(a_i) \cdot (e^\epsilon - 1) = (e^\epsilon - 1)\alpha(a) \end{aligned}$$

where the second step follows from triangle inequality, the third step follows from Eq. (3), the last step follows from definition of  $\alpha(\cdot)$  function.

Similarly, we can show

$$|\exp(a_i + b_i) - \exp(a_i)| = \exp(a_i + b_i) |\exp(-b_i) - 1| \leq \exp(a_i + b_i) \cdot (e^\epsilon - 1) \quad (4)$$

where the first step follows from simple algebra, the second step follows from Eq. (2).

Then, we have

$$\begin{aligned} |\alpha(a+b) - \alpha(a)| &= |\langle \exp(a+b), \mathbf{1}_n \rangle - \langle \exp(a), \mathbf{1}_n \rangle| \leq \sum_{i=1}^n |\exp(a_i + b_i) - \exp(a_i)| \\ &\leq \sum_{i=1}^n \exp(a_i + b_i) \cdot (e^\epsilon - 1) = (e^\epsilon - 1)\alpha(a+b) \end{aligned}$$

where the first step follows from definition of  $\alpha$  (Definition 3.1), the second step follows from triangle inequality, the third step follows from Eq. (4), and the last step follows from definition of  $\alpha$  (Definition 3.1). Thus, we complete the proof.  $\square$

## 5 RANK COLLAPSE PROPERTY

In Section 5.1, we present a Lemma which connects  $\text{Res}(\text{SAtt}())$  and  $\text{Res}()$ . In Section 5.2, we present our key lemma, a perturbation theorem for a layer of a Transformer. In Section 5.3, we present our main result and proof sketch.

### 5.1 THE CONNECTION BETWEEN $\text{Res}(\text{SAtt}())$ AND $\text{Res}()$

We next establish the relationship between  $\text{Res}(\text{SAtt}())$  and  $\text{Res}()$  in terms of the balance of the inputs.

**Lemma 5.1.** *If the following conditions hold: Let  $X \in \mathbb{R}^{n \times d}$  denote the input of attention layer. Let  $\tilde{X} = \text{SAtt}(X)$  (see Definition 1.1 for function  $\text{SAtt}$ ). Let  $W_q, W_k, W_v \in \mathbb{R}^{d \times d}$  be the weight matrices of  $\text{SAtt}$ . Let  $W = W_q W_k^\top$ . Let  $E = \beta \text{Res}(X) W \text{Res}(X)^\top$ . Suppose that  $E$  is a  $\theta$ -balanced matrix (see Definition 3.5). Let  $\beta := 1/\sqrt{d_0}$  denote the normalization factor. Let  $K := (e^\theta - 1) \|W_v\|_\infty$ . Then, we have  $\|\text{Res}(\text{SAtt}(X))\|_\infty \leq K \cdot \|\text{Res}(X)\|_\infty$ .*

*Proof.* The unscaled attention scores are computed as follows  $A = (XW_q + \mathbf{1}_n b_q^\top) \cdot (XW_k + \mathbf{1}_n b_k^\top)^\top$ . Recall that  $W = W_q W_k^\top$ . For notational convenience, we define  $b := W_k b_q$ .

We can use the softmax shift invariance property to remove terms which are constant over the columns and obtain,  $A = \underbrace{X}_{n \times d} \underbrace{W}_{d \times d} \underbrace{X^\top}_{d \times n} + \underbrace{\mathbf{1}_n}_{n \times 1} \underbrace{b^\top}_{1 \times d} \underbrace{X^\top}_{d \times n}$ .

We define  $\tilde{R} := \text{Res}(\tilde{X}) \in \mathbb{R}^{n \times d}$  (Recall the definition of the function  $\text{Res}()$  in Definition 3.4).

In next equation, we will use the definition of  $R$  to simplify  $A$ . The attention matrix can be written as

$$A = \beta \cdot (\mathbf{1}_n x^\top + R) W (\mathbf{1}_n x^\top + R)^\top + \beta \cdot \mathbf{1}_n b^\top (\mathbf{1}_n x^\top + R)^\top$$

$$= \beta \cdot (x^\top Wx \mathbf{1}_n + RWx + \mathbf{1}_n b^\top x) \mathbf{1}_n^\top + \beta \cdot (RWR^\top + \mathbf{1}_n x^\top WR^\top + \mathbf{1}_n b^\top R^\top) \quad (5)$$

Using Fact 3.2, we can remove the first term in the above equation since it is constant across columns. We thus have that the following equation for  $P = \text{softmax}(A) \in \mathbb{R}^{n \times n}$

$$P = \text{softmax}(\beta RWR^\top + \mathbf{1}_n r^\top) = \text{softmax}(E + \mathbf{1}_n r^\top) \quad (6)$$

where the first step follows from  $r = \beta R(W^\top x + b) \in \mathbb{R}^n$ , the second step follows from setting  $E = \beta RWR^\top \in \mathbb{R}^{n \times n}$ .

To continue the proof, we also set  $\tilde{A} = \mathbf{1}_n r^\top \in \mathbb{R}^{n \times n}$ , the input reweighted by the attention probabilities  $PX$  will be entry-wisely upper bounded as follows

$$\begin{aligned} PX &= P(\mathbf{1}_n x^\top + R) = \mathbf{1}_n x^\top + PR \\ &= \mathbf{1}_n x^\top + \text{softmax}(\mathbf{1}_n r^\top + E)R \\ &\leq \mathbf{1}_n x^\top + (I + e^D - I) \mathbf{1}_n \text{softmax}(r)^\top R \\ &= \mathbf{1}_n (x^\top \text{softmax}(r)^\top R) + (e^D - I) \mathbf{1}_n \text{softmax}(r)^\top R \end{aligned} \quad (7)$$

where the first step follows from definition of  $R$ , the second step follows from  $P\mathbf{1}_n = \mathbf{1}_n$ , the third step follows from Eq. (6) and  $e^D$  is a diagonal matrix, the fourth step follows from Lemma B.3 and  $e^D$  is diagonal matrix where the  $i, i$ -th entry on diagonal is  $e^{D_{i,i}}$  (see Definition 3.5 for  $D, D_{i,i} \geq 0$ ).

Therefore, the entry-wise distance of the output of the self-attention layer  $\text{SAtt}(X) = PXW_v$  from being constant across token is at most:

$$|[\text{SAtt}(X) - \mathbf{1}_n \tilde{r}^\top]_{i,j}| = |[PXW_v - \mathbf{1}_n \tilde{r}^\top]_{i,j}| \leq (e^\theta - 1) \cdot |[\mathbf{1}_n \text{softmax}(r)^\top RW_v]_{i,j}|$$

where the second step follows from  $\tilde{r} = (x + R^\top \text{softmax}(r))W_v$ , Eq. (7), and  $\theta$  (see Definition 3.5).

Now we bound the right hand side of the above inequality.

For  $\|\cdot\|_\infty$ , we can show

$$\begin{aligned} \|\text{SAtt}(X) - \mathbf{1}_n \tilde{r}^\top\|_\infty &\leq (e^\theta - 1) \|\mathbf{1}_n\|_\infty \cdot \|\text{softmax}(r)^\top RW_v\|_\infty \\ &\leq (e^\theta - 1) \|\mathbf{1}_n\|_\infty \|R\|_\infty \|W_v\|_\infty \leq (e^\theta - 1) \|R\|_\infty \cdot \|W_v\|_\infty, \end{aligned} \quad (8)$$

where the last step follows from Definition 3.5.

Note that  $R' = \text{Res}(\text{SAtt}(X))$  and  $R = \text{Res}(X)$  and using the definition of  $K$  in Lemma statement, we can show  $\|\text{Res}(\text{SAtt}(X))\|_\infty \leq K \cdot \|\text{Res}(X)\|_\infty$ . Thus, we complete the proof.  $\square$

## 5.2 PERTURBATION OF ONE TRANSFORMER LAYER

We next give a toy lemma to demonstrate the idea behind our approach. We will assume for simplicity here that  $d = n$  so  $X$  is a square matrix, and we focus on a layer with only one attention head. The full proof for any  $d \neq n$  and multiple heads appears in Lemma C.1 in the Appendix below.

**Lemma 5.2** (Single Head). *Let  $n = d$  and let  $X \in \mathbb{R}^{n \times d}$ . Let  $A = \text{softm}_1(X)$  (Recall that  $\text{softm}()$  function is defined as Definition 3.1. Note that  $\text{softm}_1$  and  $\text{softm}_2$  are two different instantiations with different  $W_k, W_q, W_v$  weights with  $\|W_k\|_\infty, \|W_q\|_\infty, \|W_v\|_\infty \leq \eta$ . Let  $B = X + A$ . Suppose  $\|\text{Res}(A)\|_\infty \leq K \cdot \|\text{Res}(X)\|_\infty \leq \epsilon$ . (We remark that this condition will hold due to Lemma 5.1; here  $K$  is as defined in Lemma 5.1). Let  $g(\epsilon) := 4\epsilon$  and let  $\epsilon_0 = 3g(3\epsilon)$ . Then we can show*

$$\|\text{softm}_2(B) - \text{softm}_2(X)\|_\infty \leq \epsilon_0.$$

*Proof.* Let  $R_X = \text{Res}(X)$  so that  $X = R_X + \mathbf{1}_n y_X^\top$  for some vector  $y_X \in \mathbb{R}^d$ . Using Lemma D.6

$$\|\text{softm}_2(X) - \text{softm}_2(R_X)\|_\infty \leq g(\epsilon) \quad (9)$$

Let  $R_A = \text{Res}(A)$  so that  $A = R_A + \mathbf{1}_n y_A^\top$  for some vector  $y_A \in \mathbb{R}^d$ . Let  $R_B = \text{Res}(B)$  so that  $B = R_B + \mathbf{1}_n y_B^\top$  for some vector  $y_B \in \mathbb{R}^d$ . Using Lemma D.6, we can show that

$$\|\text{softm}_2(B) - \text{softm}_2(R_B)\|_\infty \leq g(\epsilon). \quad (10)$$

Thus, as long as  $\|R_A\|_\infty \leq \epsilon$ , then using [Lemma D.6](#), we have

$$\|\text{softm}_2(R_X + R_A) - \text{softm}_2(R_X)\|_\infty \leq g(\epsilon) \quad (11)$$

We can show  $R_X = \text{Res}(X) = \text{Res}(B - A) = \text{Res}(B - R_A)$ . Then, we know  $\|R_X - R_B\|_\infty \leq \|\text{Res}(B - R_A) - \text{Res}(B)\|_\infty \leq 2\|R_A\|_\infty \leq 2\epsilon$ . [Here, the second step follows from Lemma 4.1.](#)

Recall  $B = X + A$ , and  $\|R_A\|_\infty \leq \epsilon$  then we know

$$\|R_X + R_A - R_B\|_\infty \leq \|R_X - R_B\|_\infty + \|R_A\|_\infty \leq 3\|R_A\|_\infty \leq 3\epsilon.$$

Since  $\|R_X + R_A - R_B\|_\infty \leq 3\epsilon$ , then using [Lemma D.6](#), we have

$$\|\text{softm}_2(R_X + R_A) - \text{softm}_2(R_B)\|_\infty \leq g(3\epsilon) \quad (12)$$

Then, we can show

$$\begin{aligned} & \|\text{softm}_2(B) - \text{softm}_2(X)\|_\infty \\ &= \|\text{softm}_2(B) - \text{softm}_2(R_X)\|_\infty + g(\epsilon) \\ &= \|\text{softm}_2(R_B) - \text{softm}_2(R_X)\|_\infty + 2g(\epsilon) \\ &\leq \|\text{softm}_2(R_B) - \text{softm}_2(R_X + R_A)\|_\infty + \|\text{softm}_2(R_X + R_A) - \text{softm}_2(R_X)\|_\infty + 2g(\epsilon) \\ &\leq g(3\epsilon) + 3g(\epsilon) \leq 3g(3\epsilon) \end{aligned}$$

where the first step follows from [triangle inequality and Eq. \(9\)](#), the second step follows from [triangle inequality and Eq. \(10\)](#), the third step follows from triangle inequality, the fourth step follows from Eq. (11) and Eq. (12), and the last step follows from  $g$  is monotone.  $\square$

### 5.3 PUTTING IT ALL TOGETHER

*Proof Sketch of Theorem 1.2.* We'll show what to do to delete one layer, then repeat that  $L - 1$  times to get down to one layer. When we delete the first layer, [Lemma C.1](#) (which is the version of [Lemma 5.2](#) which deals with multiple heads) says that the output of the second layer will differ by at most  $O(\eta \cdot \epsilon_0)$ , where  $\epsilon_0 = O(1) \cdot \|X\|_\infty$  is the constant from [Lemma 5.1](#) and [Lemma 5.2](#), and  $X$  is the input of first layer of network. Therefore, by applying [Lemma B.2](#) iteratively to each layer, it follows that the outputs of all subsequent layers will also change by at most  $O(\eta \cdot \epsilon_0)$ . In particular, the final output will differ by at most  $O(\eta \cdot \epsilon_0)$ . We finally repeat this  $L - 1$  times to remove all but one layer and get the final error. We defer further proof details to the Appendix due to space limitations.  $\square$

## 6 CONCLUSION

We have shown that Self Attention Networks must experience layer collapse unless they have large attention weights, even if they have skip connections. Our result proves that two different common notions in the literature are actually misconceptions.

The first misconception is the common interpretation of the prior work ([Dong et al., 2021](#)) that skip connections are the key to the expressive power of Self Attention Networks. We extend their result and show that even with skip connections, large weights are needed to prevent layer collapse.

The second misconception is that Self Attention Networks with smaller weights may still have reasonable expressive power. Indeed, although it is intuitive that bounding the magnitudes of weights must limit the expressive power to some extent, there is nonetheless a long line of work on trying to use networks with small weights, weight quantization, or similar approaches. This work is (presumably) hoping that the limit is only modest. We show that the limit is severe: networks with small weights cannot take advantage of more than one layer! This is the first theoretical limitation result on networks with small weights to our knowledge.

486 ETHIC STATEMENT

487  
488 This paper does not involve human subjects, personally identifiable data, or sensitive applications.  
489 We do not foresee direct ethical risks. We follow the ICLR Code of Ethics and affirm that all aspects  
490 of this research comply with the principles of fairness, transparency, and integrity.  
491

492 REPRODUCIBILITY STATEMENT

493  
494 We ensure reproducibility of our theoretical results by including all formal assumptions, definitions,  
495 and complete proofs in the appendix. The main text states each theorem clearly and refers to the  
496 detailed proofs. No external data or software is required.  
497

498 REFERENCES

- 499  
500 Kwangjun Ahn, Xiang Cheng, Minhak Song, Chulhee Yun, Ali Jadbabaie, and Suvrit Sra. Lin-  
501 ear attention is (maybe) all you need (to understand transformer optimization). In *The Twelfth*  
502 *International Conference on Learning Representations, 2024*.  
503
- 504 Josh Alman and Zhao Song. Fast attention requires bounded entries. *NeurIPS, 2023*.  
505
- 506 Josh Alman and Zhao Song. How to capture higher-order correlations? generalizing matrix softmax  
507 attention to kronecker computation. In *ICLR, 2024a*.
- 508 Josh Alman and Zhao Song. The fine-grained complexity of gradient computation for training large  
509 language models. In *NeurIPS, 2024b*.  
510
- 511 Josh Alman and Zhao Song. Fast rope attention: Combining the polynomial method and fast fourier  
512 transform. In *arXiv preprint arXiv:2505.11892, 2025*.
- 513 Hritik Bansal, Zongyu Lin, Tianyi Xie, Zeshun Zong, Michal Yarom, Yonatan Bitton, Chenfanfu  
514 Jiang, Yizhou Sun, Kai-Wei Chang, and Aditya Grover. Videophy: Evaluating physical common-  
515 sense for video generation. In *Workshop on Video-Language Models @ NeurIPS 2024, 2025*.  
516 URL <https://openreview.net/forum?id=xMlYKYFd03>.  
517
- 518 Federico Barbero, Alvaro Arroyo, Xiangming Gu, Christos Perivolaropoulos, Michael Bronstein,  
519 Razvan Pascanu, et al. Why do llms attend to the first token? *arXiv preprint arXiv:2504.02732*,  
520 2025.
- 521 Iz Beltagy, Matthew E Peters, and Arman Cohan. Longformer: The long-document transformer.  
522 *arXiv preprint arXiv:2004.05150, 2020*.  
523
- 524 Emily M Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. On the  
525 dangers of stochastic parrots: Can language models be too big? In *Proceedings of the 2021 ACM*  
526 *conference on fairness, accountability, and transparency*, pp. 610–623, 2021.
- 527 Andreas Blattmann, Tim Dockhorn, Sumith Kulal, Daniel Mendelevitch, Maciej Kilian, Dominik  
528 Lorenz, Yam Levi, Zion English, Vikram Voleti, Adam Letts, et al. Stable video diffusion: Scaling  
529 latent video diffusion models to large datasets. *arXiv preprint arXiv:2311.15127, 2023*.  
530
- 531 Matteo Bonino, Giorgia Ghione, and Giansalvo Cirrincione. The geometry of bert. *arXiv preprint*  
532 *arXiv:2502.12033, 2025*.
- 533 Christos Boutsidis and David P Woodruff. Optimal cur matrix decompositions. In *Proceedings of*  
534 *the forty-sixth annual ACM symposium on Theory of computing (STOC)*, pp. 353–362, 2014.  
535
- 536 Jan van den Brand, Binghui Peng, Zhao Song, and Omri Weinstein. Training (overparametrized)  
537 neural networks in near-linear time. *ITCS, 2021*.
- 538 Sébastien Bubeck. Convex optimization: Algorithms and complexity. *Foundations and Trends® in*  
539 *Machine Learning*, 8(3-4):231–357, 2015.

- 540 Yuefan Cao, Xuyang Guo, Jiayan Huo, Yingyu Liang, Zhenmei Shi, Zhao Song, Jiahao Zhang, and  
541 Zhen Zhuang. Text-to-image diffusion models cannot count, and prompt refinement cannot help.  
542 *arXiv preprint arXiv:2503.06884*, 2025.
- 543  
544 Thibault Castells, Hyoung-Kyu Song, Bo-Kyeong Kim, and Shinkook Choi. Ld-pruner: Efficient  
545 pruning of latent diffusion models using task-agnostic insights. In *Proceedings of the IEEE/CVF*  
546 *Conference on Computer Vision and Pattern Recognition*, pp. 821–830, 2024.
- 547 Benjamin Paul Chamberlain, Sergey Shirobokov, Emanuele Rossi, Fabrizio Frasca, Thomas  
548 Markovich, Nils Yannick Hammerla, Michael M. Bronstein, and Max Hansmire. Graph neu-  
549 ral networks for link prediction with subgraph sketching. In *The Eleventh International Confer-*  
550 *ence on Learning Representations*, 2023. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=mloqEOAozQU)  
551 [mloqEOAozQU](https://openreview.net/forum?id=mloqEOAozQU).
- 552 Ya-Ting Chang, Zhibo Hu, Xiaoyu Li, Shuiqiao Yang, Jiaojiao Jiang, and Nan Sun. Dihan: A novel  
553 dynamic hierarchical graph attention network for fake news detection. In *Proceedings of the 33rd*  
554 *ACM International Conference on Information and Knowledge Management*, pp. 197–206, 2024.
- 555  
556 Moses Charikar, Michael Kapralov, Navid Nouri, and Paris Siminelakis. Kernel density estimation  
557 through density constrained near neighbor search. In *2020 IEEE 61st Annual Symposium on*  
558 *Foundations of Computer Science (FOCS)*, pp. 172–183. IEEE, 2020.
- 559 Bo Chen, Xiaoyu Li, Yingyu Liang, Jiangxuan Long, Zhenmei Shi, and Zhao Song. Circuit com-  
560 plexity bounds for rope-based transformer architecture. *arXiv preprint arXiv:2411.07602*, 2024a.
- 561  
562 Bo Chen, Yingyu Liang, Zhizhou Sha, Zhenmei Shi, and Zhao Song. Hsr-enhanced sparse attention  
563 acceleration. *arXiv preprint arXiv:2410.10165*, 2024b.
- 564 Bo Chen, Chengyue Gong, Xiaoyu Li, Yingyu Liang, Zhizhou Sha, Zhenmei Shi, Zhao Song,  
565 and Mingda Wan. High-order matching for one-step shortcut diffusion models. *arXiv preprint*  
566 *arXiv:2502.00688*, 2025a.
- 567  
568 Jingye Chen, Yupan Huang, Tengchao Lv, Lei Cui, Qifeng Chen, and Furu Wei. Textdiffuser:  
569 Diffusion models as text painters. *Advances in Neural Information Processing Systems*, 36:9353–  
570 9387, 2023.
- 571  
572 Lei Chen, Yuan Meng, Chen Tang, Xinzhu Ma, Jingyan Jiang, Xin Wang, Zhi Wang, and Wenwu  
573 Zhu. Q-dit: Accurate post-training quantization for diffusion transformers. In *Proceedings of the*  
*Computer Vision and Pattern Recognition Conference (CVPR)*, pp. 28306–28315, 2025b.
- 574  
575 Yifang Chen, Jiayan Huo, Xiaoyu Li, Yingyu Liang, Zhenmei Shi, and Zhao Song. Fast gradient  
576 computation for rope attention in almost linear time. *arXiv preprint arXiv:2412.17316*, 2024c.
- 577  
578 Eli Chien, Chao Pan, and Olgica Milenkovic. Efficient model updates for approximate unlearning  
579 of graph-structured data. In *The Eleventh International Conference on Learning Representations*,  
2023. URL <https://openreview.net/forum?id=fhcu4FBLciL>.
- 580  
581 Rewon Child, Scott Gray, Alec Radford, and Ilya Sutskever. Generating long sequences with sparse  
582 transformers. *arXiv preprint arXiv:1904.10509*, 2019.
- 583  
584 Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi  
585 Wang, Mostafa Dehghani, Siddhartha Brahma, et al. Scaling instruction-finetuned language mod-  
els. *arXiv preprint arXiv:2210.11416*, 2022.
- 586  
587 Kenneth L Clarkson and David P Woodruff. Low-rank approximation and regression in input spar-  
sity time. In *STOC*, 2013.
- 588  
589 Enyan Dai, Wei Jin, Hui Liu, and Suhang Wang. Towards robust graph neural networks for noisy  
590 graphs with sparse labels. In *Proceedings of the fifteenth ACM international conference on web*  
591 *search and data mining*, pp. 181–191, 2022.
- 592  
593 Trung Dao, Thuan Hoang Nguyen, Thanh Le, Duc Vu, Khoi Nguyen, Cuong Pham, and Anh Tran.  
Swiftbrush v2: Make your one-step diffusion model better than its teacher. In *European Confer-*  
*ence on Computer Vision*, pp. 176–192. Springer, 2024.

- 594 Juncan Deng, Shuaiting Li, Zeyu Wang, Hong Gu, Kedong Xu, and Kejie Huang. Vq4dit: Effi-  
595 cient post-training vector quantization for diffusion transformers. In *Proceedings of the AAAI*  
596 *Conference on Artificial Intelligence (AAAI)*, volume 39:15, pp. 16226–16234, 2025.
- 597 Yichuan Deng, Zhao Song, Zifan Wang, and Han Zhang. Streaming kernel pca algorithm with small  
598 space. *arXiv preprint arXiv:2303.04555*, 2023a.
- 600 Yichuan Deng, Zhao Song, and Junze Yin. Faster robust tensor power method for arbitrary order.  
601 *arXiv preprint arXiv:2306.00406*, 2023b.
- 602 Yichuan Deng, Zhao Song, and Tianyi Zhou. Superiority of softmax: Unveiling the performance  
603 edge over linear attention. *arXiv preprint arXiv:2310.11685*, 2023c.
- 604 Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. Qlora: Efficient finetuning  
605 of quantized llms. *Advances in neural information processing systems (NeurIPS)*, 36:10088–  
606 10115, 2023.
- 607 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep  
608 bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- 609 Huaian Diao, Zhao Song, Wen Sun, and David Woodruff. Sketching for kronecker product re-  
610 gression and p-splines. In *International Conference on Artificial Intelligence and Statistics*, pp.  
611 1299–1308. PMLR, 2018.
- 612 Huaian Diao, Rajesh Jayaram, Zhao Song, Wen Sun, and David Woodruff. Optimal sketching  
613 for kronecker product regression and low rank approximation. *Advances in neural information*  
614 *processing systems*, 32, 2019.
- 615 Mucong Ding, Kezhi Kong, Jingling Li, Chen Zhu, John Dickerson, Furong Huang, and Tom Gold-  
616 stein. Vq-gnn: A universal framework to scale up graph neural networks using vector quantiza-  
617 tion. *Advances in Neural Information Processing Systems*, 34:6733–6746, 2021.
- 618 Mucong Ding, Tahseen Rabbani, Bang An, Evan Wang, and Furong Huang. Sketch-gnn: Scalable  
619 graph neural networks with sublinear training complexity. In *Advances in Neural Information*  
620 *Processing Systems*, 2022.
- 621 Yihe Dong, Jean-Baptiste Cordonnier, and Andreas Loukas. Attention is not all you need: Pure  
622 attention loses rank doubly exponentially with depth. In *International conference on machine*  
623 *learning*, pp. 2793–2803. PMLR, 2021.
- 624 Kazuki Egashira, Mark Vero, Robin Staab, Jingxuan He, and Martin Vechev. Exploiting llm quanti-  
625 zation. *Advances in Neural Information Processing Systems (NeurIPS)*, 37:41709–41732, 2024.
- 626 Wenqi Fan, Yao Ma, Qing Li, Yuan He, Eric Zhao, Jiliang Tang, and Dawei Yin. Graph neural  
627 networks for social recommendation. In *The world wide web conference*, pp. 417–426, 2019.
- 628 Kevin Frans, Danijar Hafner, Sergey Levine, and Pieter Abbeel. One step diffusion via shortcut  
629 models. *arXiv preprint arXiv:2410.12557*, 2024.
- 630 Elias Frantar and Dan Alistarh. Sparsegpt: Massive language models can be accurately pruned in  
631 one-shot. In *International Conference on Machine Learning*, pp. 10323–10337. PMLR, 2023.
- 632 Elias Frantar, Saleh Ashkboos, Torsten Hoefler, and Dan Alistarh. Gptq: Accurate post-training  
633 quantization for generative pre-trained transformers. *arXiv preprint arXiv:2210.17323*, 2022.
- 634 Minghao Fu, Hao Yu, Jie Shao, Junjie Zhou, Ke Zhu, and Jianxin Wu. Quantization without tears. In  
635 *Proceedings of the Computer Vision and Pattern Recognition Conference*, pp. 4462–4472, 2025.
- 636 Ziwang Fu, Feng Liu, Jiahao Zhang, Hanyang Wang, Chengyi Yang, Qing Xu, Jiayin Qi, Xiangling  
637 Fu, and Aimin Zhou. Sagn: semantic adaptive graph network for skeleton-based human action  
638 recognition. In *Proceedings of the 2021 International Conference on Multimedia Retrieval*, pp.  
639 110–117, 2021.

- 648 Peng Gao, Jiaming Han, Renrui Zhang, Ziyi Lin, Shijie Geng, Aojun Zhou, Wei Zhang, Pan Lu,  
649 Conghui He, Xiangyu Yue, et al. Llama-adapter v2: Parameter-efficient visual instruction model.  
650 *arXiv preprint arXiv:2304.15010*, 2023a.
- 651 Tianyu Gao, Adam Fisch, and Danqi Chen. Making pre-trained language models better few-shot  
652 learners. In *Proceedings of the 59th Annual Meeting of the Association for Computational Lin-*  
653 *guistics and the 11th International Joint Conference on Natural Language Processing*, 2021.
- 654 Yeqi Gao, Sridhar Mahadevan, and Zhao Song. An over-parameterized exponential regression.  
655 *arXiv preprint arXiv:2303.16504*, 2023b.
- 656 Yeqi Gao, Zhao Song, Xin Yang, and Ruizhe Zhang. Fast quantum algorithm for attention compu-  
657 tation. *arXiv preprint arXiv:2307.08045*, 2023c.
- 658 Yeqi Gao, Zhao Song, Xin Yang, and Yufa Zhou. Differentially private attention computation. In  
659 *Neurips Safe Generative AI Workshop 2024*, 2024.
- 660 Simon Geisler, Tobias Schmidt, Hakan Sirin, Daniel Zügner, Aleksandar Bojchevski, and Stephan  
661 Günnemann. Robustness of graph neural networks at scale. In *NeurIPS*, 2021.
- 662 Borjan Geshkovski, Cyril Letrouit, Yury Polyanskiy, and Philippe Rigollet. The emergence of clus-  
663 ters in self-attention dynamics. *Advances in Neural Information Processing Systems*, 36:57026–  
664 57037, 2023.
- 665 Chengyue Gong, Yekun Ke, Xiaoyu Li, Yingyu Liang, Zhizhou Sha, Zhenmei Shi, and Zhao  
666 Song. On computational limits of flowar models: Expressivity and efficiency. *arXiv preprint*  
667 *arXiv:2502.16490*, 2025.
- 668 Google. Gemini 2.0 is now available to everyone, 2025. URL [https://blog.google/  
669 technology/google-deepmind/gemini-model-updates-february-2025/](https://blog.google/technology/google-deepmind/gemini-model-updates-february-2025/).
- 670 Yuzhou Gu, Zhao Song, Junze Yin, and Lichen Zhang. Low rank matrix completion via robust al-  
671 ternating minimization in nearly linear time. In *Proceedings of the 12th International Conference  
672 on Learning Representations, ICLR'24*, 2024.
- 673 Yuzhou Gu, Zhao Song, and Lichen Zhang. Faster algorithms for structured linear and kernel support  
674 vector machines. In *The Thirteenth International Conference on Learning Representations*, 2025.  
675 URL <https://openreview.net/forum?id=DDNFTaVQdU>.
- 676 Xiaojun Guo, Yifei Wang, Tianqi Du, and Yisen Wang. Contranorm: A contrastive learning per-  
677 spective on oversmoothing and beyond. *arXiv preprint arXiv:2303.06562*, 2023.
- 678 Xuyang Guo, Zekai Huang, Jiayan Huo, Yingyu Liang, Zhenmei Shi, Zhao Song, and Jiahao Zhang.  
679 Can you count to nine? a human evaluation benchmark for counting limits in modern text-to-video  
680 models. *arXiv preprint arXiv:2504.04051*, 2025a.
- 681 Xuyang Guo, Jiayan Huo, Zhenmei Shi, Zhao Song, Jiahao Zhang, and Jiale Zhao. T2vphysbench:  
682 A first-principles benchmark for physical consistency in text-to-video generation. *arXiv preprint*  
683 *arXiv:2505.00337*, 2025b.
- 684 Xuyang Guo, Jiayan Huo, Zhenmei Shi, Zhao Song, Jiahao Zhang, and Jiale Zhao. T2vtextbench:  
685 A human evaluation benchmark for textual control in video generation models. *arXiv preprint*  
686 *arXiv:2505.04946*, 2025c.
- 687 Shreya Gupta, Boyang Huang, Barna Saha, Yinzhan Xu, and Christopher Ye. Subquadratic algo-  
688 rithms and hardness for attention with any temperature. In *arXiv preprint arXiv:2505.14840*,  
689 2025.
- 690 Will Hamilton, Zhitao Ying, and Jure Leskovec. Inductive representation learning on large graphs.  
691 *Advances in neural information processing systems*, 30, 2017.
- 692 Insu Han, Rajesh Jayaram, Amin Karbasi, Vahab Mirrokni, David Woodruff, and Amir Zandieh.  
693 Hyperattention: Long-context attention in near-linear time. In *The Twelfth International Confer-  
694 ence on Learning Representations*, 2024.

- 702 DongNyeong Heo and Heeyoul Choi. Generalized probabilistic attention mechanism in transform-  
703 ers. *arXiv preprint arXiv:2410.15578*, 2024.  
704
- 705 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in*  
706 *neural information processing systems*, 33:6840–6851, 2020.  
707
- 708 Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yanzhi Li, Shean Wang, Lu Wang,  
709 and Weizhu Chen. LoRA: Low-rank adaptation of large language models. In *International Con-*  
710 *ference on Learning Representations*, 2022.  
711
- 712 Jerry Yao-Chieh Hu, Thomas Lin, Zhao Song, and Han Liu. On computational limits of modern  
713 hopfield models: A fine-grained complexity analysis. In *Forty-first International Conference on*  
714 *Machine Learning (ICML)*, 2024a.  
715
- 716 Jerry Yao-Chieh Hu, Maojiang Su, En-Jui Kuo, Zhao Song, and Han Liu. Computational limits  
717 of low-rank adaptation (lora) for transformer-based models. *arXiv preprint arXiv:2406.03136*,  
718 2024b.  
719
- 720 Jerry Yao-Chieh Hu, Dennis Wu, and Han Liu. Provably optimal memory capacity for modern hop-  
721 field models: Tight analysis for transformer-compatible dense associative memories. In *Advances*  
722 *in Neural Information Processing Systems (NeurIPS)*, volume 37, 2024c.  
723
- 724 Jerry Yao-Chieh Hu, Weimin Wu, Zhao Song, and Han Liu. On statistical rates and provably efficient  
725 criteria of latent diffusion transformers (dits). *arXiv preprint arXiv:2407.01079*, 2024d.  
726
- 727 Jerry Yao-Chieh Hu, Weimin Wu, Yi-Chen Lee, Yu-Chao Huang, Minshuo Chen, and Han Liu. On  
728 statistical rates of conditional diffusion transformers: Approximation, estimation and minimax  
729 optimality. In *The Thirteenth International Conference on Learning Representations*, 2025a.  
730
- 731 Weiming Hu, Haoyan Zhang, Cong Guo, Yu Feng, Renyang Guan, Zhendong Hua, Zihan Liu,  
732 Yue Guan, Minyi Guo, and Jingwen Leng. M-ant: Efficient low-bit group quantization for llms  
733 via mathematically adaptive numerical type. In *2025 IEEE International Symposium on High*  
734 *Performance Computer Architecture (HPCA)*, pp. 1112–1126. IEEE, 2025b.  
735
- 736 Xing Hu, Yuan Cheng, Dawei Yang, Zhixuan Chen, Zukang Xu, Jiangyong Yu, XUCHEN, Zhi-  
737 hang Yuan, Zhe jiang, and Sifan Zhou. OSTQuant: Refining large language model quanti-  
738 zation with orthogonal and scaling transformations for better distribution fitting. In *The Thir-*  
739 *teenth International Conference on Learning Representations (ICLR)*, 2025c. URL <https://openreview.net/forum?id=rAcgDBdKnP>.  
740
- 741 Itay Hubara, Brian Chmiel, Moshe Island, Ron Banner, Joseph Naor, and Daniel Soudry. Accelerated  
742 sparse neural training: A provable and efficient method to find n: m transposable masks.  
743 *Advances in neural information processing systems*, 34:21099–21111, 2021.  
744
- 745 Xiaofei Hui, Qian Wu, Hossein Rahmani, and Jun Liu. Class-agnostic object counting with text-to-  
746 image diffusion model. In *European Conference on Computer Vision*, pp. 1–18. Springer, 2024.  
747
- 748 Yiping Ji, Hemanth Saratchandran, Peyman Moghaddam, and Simon Lucey. Always skip attention.  
749 *arXiv preprint arXiv:2505.01996*, 2025.  
750
- 751 Shunhua Jiang, Zhao Song, Omri Weinstein, and Hengjie Zhang. A faster algorithm for solving gen-  
752 eral lps. In *Proceedings of the 53rd Annual ACM SIGACT Symposium on Theory of Computing*,  
753 pp. 823–832, 2021.  
754
- 755 Angelos Katharopoulos, Apoorv Vyas, Nikolaos Pappas, and François Fleuret. Transformers are  
rnn: Fast autoregressive transformers with linear attention. In *International conference on ma-*  
*chine learning*, pp. 5156–5165. PMLR, 2020a.
- Angelos Katharopoulos, Apoorv Vyas, Nikolaos Pappas, and François Fleuret. Transformers are  
rnn: Fast autoregressive transformers with linear attention. In *International conference on ma-*  
*chine learning*, pp. 5156–5165. PMLR, 2020b.

- 756 Jonathan A Kelner, Frederic Koehler, Raghu Meka, and Dhruv Rohatgi. On the power of precon-  
757 ditioning in sparse linear regression. In *2021 IEEE 62nd Annual Symposium on Foundations of*  
758 *Computer Science (FOCS)*, pp. 550–561. IEEE, 2022.
- 759 Omar Khattab, Keshav Santhanam, Xiang Lisa Li, David Hall, Percy Liang, Christopher Potts,  
760 and Matei Zaharia. Demonstrate-search-predict: Composing retrieval and language models for  
761 knowledge-intensive nlp. *arXiv preprint arXiv:2212.14024*, 2022.
- 762 Bo-Kyeong Kim, Geonmin Kim, Tae-Ho Kim, Thibault Castells, Shinkook Choi, Junho Shin, and  
763 Hyoung-Kyu Song. Shortened llama: A simple depth pruning for large language models. *arXiv*  
764 *preprint arXiv:2402.02834*, 11, 2024.
- 765 Jihwan Kim, Miso Lee, and Jae-Pil Heo. Self-feedback detr for temporal action detection. In  
766 *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 10286–10296,  
767 2023.
- 768 Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional net-  
769 works. *ICLR*, 2017.
- 770 Eldar Kurtic, Denis Kuznedelev, Elias Frantar, Michael Goin, and Dan Alistarh. Sparse finetuning  
771 for inference acceleration of large language models. *arXiv preprint arXiv:2310.06927*, 2023.
- 772 Chenyang Li, Yingyu Liang, Zhenmei Shi, and Zhao Song. Exploring the frontiers of softmax: Pro-  
773 vable optimization, applications in diffusion model, and beyond. *arXiv preprint arXiv:2405.03251*,  
774 2024a.
- 775 Chenyang Li, Yingyu Liang, Zhenmei Shi, and Zhao Song. When can we solve the weighted low  
776 rank approximation problem in truly subquadratic time? In *AISTATS*, 2025a.
- 777 Maosen Li, Siheng Chen, Xu Chen, Ya Zhang, Yanfeng Wang, and Qi Tian. Symbiotic graph  
778 neural networks for 3d skeleton-based human action recognition and motion prediction. *IEEE*  
779 *transactions on pattern analysis and machine intelligence*, 44(6):3316–3333, 2021.
- 780 Miaoyu Li, Ying Fu, and Yulun Zhang. Spatial-spectral transformer for hyperspectral image de-  
781 noising. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37:1, pp.  
782 1368–1376, 2023.
- 783 Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation.  
784 In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics*  
785 *and the 11th International Joint Conference on Natural Language Processing*. Association for  
786 Computational Linguistics, 2021.
- 787 Xiaoyu Li, Yingyu Liang, Zhenmei Shi, and Zhao Song. A tighter complexity analysis of sparsegpt.  
788 *arXiv preprint arXiv:2408.12151*, 2024b.
- 789 Xiaoyu Li, Yingyu Liang, Jiangxuan Long, Zhenmei Shi, Zhao Song, and Zhen Zhuang. Simulation  
790 of hypergraph algorithms with looped transformers. *arXiv preprint arXiv:2501.10688*, 2025b.
- 791 Yuanzhi Li, Yingyu Liang, and Andrej Risteski. Recovery guarantee of weighted low-rank ap-  
792 proximation via alternating minimization. In *International Conference on Machine Learning*, pp.  
793 2358–2367. PMLR, 2016.
- 794 Yingyu Liang, Heshan Liu, Zhenmei Shi, Zhao Song, and Junze Yin. Conv-basis: A new  
795 paradigm for efficient attention inference and gradient computation in transformers. *arXiv*  
796 *preprint arXiv:2405.05219*, 2024a.
- 797 Yingyu Liang, Zhizhou Sha, Zhenmei Shi, and Zhao Song. Differential privacy mechanisms in  
798 neural tangent kernel regression. *arXiv preprint arXiv:2407.13621*, 2024b.
- 799 Yingyu Liang, Zhizhou Sha, Zhenmei Shi, Zhao Song, and Yufa Zhou. Multi-layer transformers  
800 gradient can be approximated in almost linear time. *arXiv preprint arXiv:2408.13233*, 2024c.
- 801 Yingyu Liang, Zhenmei Shi, Zhao Song, and Yufa Zhou. Differential privacy of cross-attention with  
802 provable guarantee. *arXiv preprint arXiv:2407.14717*, 2024d.

- 810 Yingyu Liang, Zhenmei Shi, Zhao Song, and Yufa Zhou. Tensor attention training: Provably effi-  
811 cient learning of higher-order transformers. *arXiv preprint arXiv:2405.16411*, 2024e.  
812
- 813 Yingyu Liang, Jiangxuan Long, Zhenmei Shi, Zhao Song, and Yufa Zhou. Beyond linear approxima-  
814 tions: A novel pruning approach for attention matrix. In *The Thirteenth International Conference*  
815 *on Learning Representations*, 2025.
- 816 Ji Lin, Jiaming Tang, Haotian Tang, Shang Yang, Guangxuan Xiao, and Song Han. Awq: Activation-  
817 aware weight quantization for on-device llm compression and acceleration. *GetMobile: Mobile*  
818 *Computing and Communications*, 28(4):12–17, 2025.  
819
- 820 Wanyu Lin, Baochun Li, and Cong Wang. Towards private learning on decentralized graphs with  
821 local differential privacy. *IEEE Transactions on Information Forensics and Security*, 17:2936–  
822 2946, 2022.
- 823 Zhenghao Lin, Yeyun Gong, Yelong Shen, Tong Wu, Zhihao Fan, Chen Lin, Nan Duan, and Weizhu  
824 Chen. Text generation with diffusion language models: A pre-training approach with continuous  
825 paragraph denoise. In *International Conference on Machine Learning*, pp. 21051–21064. PMLR,  
826 2023.
- 827 Chengyi Liu, Jiahao Zhang, Shijie Wang, Wenqi Fan, and Qing Li. Score-based generative diffusion  
828 models for social recommendations. *arXiv preprint arXiv:2412.15579*, 2024a.  
829
- 830 Han Liu, Haotian Gao, Xiaotong Zhang, Changya Li, Feng Zhang, Wei Wang, Fenglong Ma, and  
831 Hong Yu. Septq: A simple and effective post-training quantization paradigm for large language  
832 models. In *Proceedings of the 31st ACM SIGKDD Conference on Knowledge Discovery and Data*  
833 *Mining V. 1 (KDD)*, pp. 812–823, 2025.
- 834 Tianlin Liu and Friedemann Zenke. Finding trainable sparse networks through neural tangent trans-  
835 fer. In *International Conference on Machine Learning*, pp. 6336–6347. PMLR, 2020.  
836
- 837 Zechun Liu, Changsheng Zhao, Igor Fedorov, Bilge Soran, Dhruv Choudhary, Raghuraman Krish-  
838 namoorthi, Vikas Chandra, Yuandong Tian, and Tijmen Blankevoort. Spinqant: Llm quantiza-  
839 tion with learned rotations. *arXiv preprint arXiv:2405.16406*, 2024b.
- 840 AI @ Meta Llama Team. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.  
841
- 842 Wenchi Ma, Tianxiao Zhang, and Guanghui Wang. Miti-detr: Object detection based on transform-  
843 ers with mitigatory self-attention convergence. *arXiv preprint arXiv:2112.13310*, 2021.  
844
- 845 Xinyin Ma, Gongfan Fang, and Xinchao Wang. Deepcache: Accelerating diffusion models for  
846 free. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*,  
847 pp. 15762–15772, 2024.
- 848 Joseph McDonald, Baolin Li, Nathan Frey, Devesh Tiwari, Vijay Gadepally, and Siddharth Samsi.  
849 Great power, great responsibility: Recommendations for reducing energy for training language  
850 models. In *Findings of the Association for Computational Linguistics: NAACL 2022*, pp. 1962–  
851 1970, 2022.
- 852 Swaroop Mishra, Daniel Khashabi, Chitta Baral, and Hannaneh Hajishirzi. Cross-task generalization  
853 via natural language crowdsourcing instructions. In *Proceedings of the 60th Annual Meeting of*  
854 *the Association for Computational Linguistics*, 2022.  
855
- 856 Saman Motamed, Laura Culp, Kevin Swersky, Priyank Jaini, and Robert Geirhos. Do generative  
857 video models learn physical principles from watching videos? *arXiv preprint arXiv:2501.09038*,  
858 2025.
- 859 Tamara T Mueller, Johannes C Paetzold, Chinmay Prabhakar, Dmitrii Usynin, Daniel Rueckert, and  
860 Georgios Kaissis. Differentially private graph neural networks for whole-graph classification.  
861 *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(6):7308–7318, 2022.  
862
- 863 Alireza Naderi, Thiziri Nait Saada, and Jared Tanner. Mind the gap: a spectral analysis of rank  
collapse and signal propagation in transformers. *arXiv preprint arXiv:2410.07799*, 2024.

- 864 Jelani Nelson and Huy L Nguyễn. Osnap: Faster numerical linear algebra algorithms via sparser  
865 subspace embeddings. In *2013 IEEE 54th Annual Symposium on Foundations of Computer Science*,  
866 pp. 117–126. IEEE, 2013.
- 867 Tam Minh Nguyen, César A Uribe, Tan Minh Nguyen, and Richard Baraniuk. Pidformer: Trans-  
868 former meets control theory. In *Forty-first International Conference on Machine Learning*, 2024.
- 869 Yotam Nitzan, Zongze Wu, Richard Zhang, Eli Shechtman, Daniel Cohen-Or, Taesung Park, and  
870 Michaël Gharbi. Lazy diffusion transformer for interactive image editing. In *European Confer-  
871 ence on Computer Vision*, pp. 55–72. Springer, 2024.
- 872 Lorenzo Noci, Sotiris Anagnostidis, Luca Biggio, Antonio Orvieto, Sidak Pal Singh, and Aurelien  
873 Lucchi. Signal propagation in transformers: Theoretical perspectives and the role of rank collapse.  
874 *Advances in Neural Information Processing Systems*, 35:27198–27211, 2022.
- 875 OpenAI. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- 876 Von Johannes Oswald, Eyvind Niklasson, Ettore Randazzo, João Sacramento, Alexander Mordv-  
877 intsev, Andrey Zhmoginov, and Max Vladymyrov. Transformers learn in-context by gradient  
878 descent. In *International Conference on Machine Learning*. PMLR, 2023.
- 879 Xu Ouyang, Tao Ge, Thomas Hartvigsen, Zhisong Zhang, Haitao Mi, and Dong Yu. Low-bit quan-  
880 tization favors undertrained llms. In *Proceedings of the 63rd Annual Meeting of the Association  
881 for Computational Linguistics (ACL)*, pp. 32338–32348, 2025.
- 882 Yeonhong Park, Jake Hyun, Hojoon Kim, and Jae W Lee. {DecDEC}: A systems approach to  
883 advancing {Low-Bit}{LLM} quantization. In *19th USENIX Symposium on Operating Systems  
884 Design and Implementation (OSDI 25)*, pp. 803–819, 2025.
- 885 Wei Peng, Xiaopeng Hong, Haoyu Chen, and Guoying Zhao. Learning graph convolutional network  
886 for skeleton-based human action recognition by neural searching. In *Proceedings of the AAAI  
887 conference on artificial intelligence*, 2020.
- 888 Sergio P Perez, Yan Zhang, James Briggs, Charlie Blake, Josh Levy-Kramer, Paul Balanca, Carlo  
889 Luschi, Stephen Barlow, and Andrew William Fitzgibbon. Training and inference of large lan-  
890 guage models using 8-bit floating point. *arXiv preprint arXiv:2309.17224*, 2023.
- 891 Ilya Razenshteyn, Zhao Song, and David P Woodruff. Weighted low rank approximations with  
892 provable guarantees. In *Proceedings of the forty-eighth annual ACM symposium on Theory of  
893 Computing*, pp. 250–263, 2016.
- 894 Aravind Reddy, Zhao Song, and Lichen Zhang. Dynamic tensor product regression. In *Conference  
895 on Neural Information Processing Systems (NeurIPS)*, pp. 4791–4804, 2022.
- 896 Andreas Roth and Thomas Liebig. Rank collapse causes over-smoothing and over-correlation in  
897 graph neural networks. In *Learning on Graphs Conference*, pp. 35–1. PMLR, 2024.
- 898 Subham Sahoo, Marianne Arriola, Yair Schiff, Aaron Gokaslan, Edgar Marroquin, Justin Chiu,  
899 Alexander Rush, and Volodymyr Kuleshov. Simple and effective masked diffusion language  
900 models. *Advances in Neural Information Processing Systems*, 37:130136–130184, 2024.
- 901 Michael E Sander, Pierre Ablin, Mathieu Blondel, and Gabriel Peyré. Sinkformers: Transform-  
902 ers with doubly stochastic attention. In *International Conference on Artificial Intelligence and  
903 Statistics*, pp. 3515–3530. PMLR, 2022.
- 904 Imanol Schlag, Kazuki Irie, and Jürgen Schmidhuber. Linear transformers are secretly fast weight  
905 programmers. In *International Conference on Machine Learning*. PMLR, 2021.
- 906 Xuan Shen, Zhao Song, Yufa Zhou, Bo Chen, Jing Liu, Ruiyi Zhang, Ryan A Rossi, Hao Tan, Tong  
907 Yu, Xiang Chen, et al. Numerical pruning for efficient autoregressive models. *arXiv preprint  
908 arXiv:2412.12441*, 2024.
- 909 Xuan Shen, Zhao Song, Yufa Zhou, Bo Chen, Yanyu Li, Yifan Gong, Kai Zhang, Hao Tan, Jason  
910 Kuen, Henghui Ding, et al. Lazydit: Lazy learning for the acceleration of diffusion transformers.  
911 In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2025.

- 918 Zhenmei Shi, Jiefeng Chen, Kunyang Li, Jayaram Raghuram, Xi Wu, Yingyu Liang, and Somesh  
919 Jha. The trade-off between universality and label efficiency of representations from contrastive  
920 learning. In *The Eleventh International Conference on Learning Representations*, 2023a.  
921
- 922 Zhenmei Shi, Jiefeng Chen, Kunyang Li, Jayaram Raghuram, Xi Wu, Yingyu Liang, and Somesh  
923 Jha. The trade-off between universality and label efficiency of representations from contrastive  
924 learning. In *The Eleventh International Conference on Learning Representations*, 2023b.
- 925 Zhenmei Shi, Junyi Wei, Zhuoyan Xu, and Yingyu Liang. Why larger language models do in-  
926 context learning differently? In *R0-FoMo: Robustness of Few-shot and Zero-shot Learning in  
927 Large Foundation Models*, 2023c.  
928
- 929 Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben  
930 Poole. Score-based generative modeling through stochastic differential equations. In *Internat-  
931 ional Conference on Learning Representations*, 2021a. URL [https://openreview.net/  
932 forum?id=PXTIG12RRHS](https://openreview.net/forum?id=PXTIG12RRHS).
- 933 Zhao Song and Zheng Yu. Oblivious sketching-based central path method for linear programming.  
934 In *International Conference on Machine Learning*, pp. 9835–9847. PMLR, 2021.  
935
- 936 Zhao Song, David P Woodruff, and Peilin Zhong. Low rank approximation with entrywise  $\ell_1$ -norm  
937 error. In *Proceedings of the 49th Annual ACM SIGACT Symposium on Theory of Computing*, pp.  
938 688–701, 2017.
- 939 Zhao Song, David Woodruff, and Peilin Zhong. Towards a zero-one law for column subset selection.  
940 *Advances in Neural Information Processing Systems*, 32, 2019a.  
941
- 942 Zhao Song, David Woodruff, and Peilin Zhong. Average case column subset selection for entrywise  
943  $\ell_1$ -norm loss. *Advances in Neural Information Processing Systems*, 32, 2019b.  
944
- 945 Zhao Song, David P Woodruff, and Peilin Zhong. Relative error tensor low rank approximation. In  
946 *Proceedings of the Thirtieth Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pp.  
947 2772–2789. SIAM, 2019c.
- 948 Zhao Song, David Woodruff, Zheng Yu, and Lichen Zhang. Fast sketching of polynomial kernels of  
949 polynomial degree. In *International Conference on Machine Learning*, pp. 9812–9823. PMLR,  
950 2021b.  
951
- 952 Zhao Song, Mingquan Ye, Junze Yin, and Lichen Zhang. A nearly-optimal bound for fast regression  
953 with  $\ell_\infty$  guarantee. In *International Conference on Machine Learning*, pp. 32463–32482. PMLR,  
954 2023a.
- 955 Zhao Song, Junze Yin, and Ruizhe Zhang. Revisiting quantum algorithms for linear regres-  
956 sions: Quadratic speedups without data-dependent parameters. *arXiv preprint arXiv:2311.14823*,  
957 2023b.  
958
- 959 Zhao Song, Junze Yin, and Lichen Zhang. Solving attention kernel regression problem via pre-  
960 conditioner. In *International Conference on Artificial Intelligence and Statistics (AISTATS)*, pp.  
961 208–216. PMLR, 2024a.
- 962 Zhao Song, Lichen Zhang, and Ruizhe Zhang. Training multi-layer over-parametrized neural net-  
963 work in subquadratic time. *ITCS*, 2024b.  
964
- 965 Zhao Song, Mingquan Ye, Junze Yin, and Lichen Zhang. Efficient alternating minimization with  
966 applications to weighted low rank approximation. In *The Thirteenth International Confer-  
967 ence on Learning Representations*, 2025. URL [https://openreview.net/forum?id=  
968 rvhu4V7yrX](https://openreview.net/forum?id=rvhu4V7yrX).
- 969
- 970 Niranjan Subrahmanya and Yung C Shin. Sparse multiple kernel learning for signal processing  
971 applications. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 32(5):788–798,  
2009.

- 972 Mingjie Sun, Zhuang Liu, Anna Bair, and J Zico Kolter. A simple and effective pruning approach  
973 for large language models. In *The Twelfth International Conference on Learning Representations*,  
974 2024.
- 975 Tianxiang Sun, Yunfan Shao, Hong Qian, Xuanjing Huang, and Xipeng Qiu. Black-box tuning for  
976 language-model-as-a-service. In *International Conference on Machine Learning*. PMLR, 2022.
- 977  
978 Yutao Sun, Li Dong, Shaohan Huang, Shuming Ma, Yuqing Xia, Jilong Xue, Jianyong Wang, and  
979 Furu Wei. Retentive network: A successor to transformer for large language models. *arXiv*  
980 *preprint arXiv:2307.08621*, 2023.
- 981  
982 Yao-Hung Hubert Tsai, Shaojie Bai, Makoto Yamada, Louis-Philippe Morency, and Ruslan  
983 Salakhutdinov. Transformer dissection: a unified understanding of transformer’s attention via  
984 the lens of kernel. *arXiv preprint arXiv:1908.11775*, 2019.
- 985  
986 Yuxiang Tuo, Wangmeng Xiang, Jun-Yan He, Yifeng Geng, and Xuansong Xie. Anytext: Multi-  
987 lingual visual text generation and editing. In *The Twelfth International Conference on Learning*  
988 *Representations*, 2024. URL <https://openreview.net/forum?id=ezBH9WE9s2>.
- 989  
990 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,  
991 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural informa-*  
*tion processing systems*, 30, 2017.
- 992  
993 Swagath Venkataramani, Anand Raghunathan, Jie Liu, and Mohammed Shoaib. Scalable-effort  
994 classifiers for energy-efficient machine learning. In *Proceedings of the 52nd annual design au-*  
*tomation conference*, pp. 1–6, 2015.
- 995  
996 Limei Wang, Kaveh Hassani, Si Zhang, Dongqi Fu, Baichuan Yuan, Weilin Cong, Zhigang Hua,  
997 Hao Wu, Ning Yao, and Bo Long. Learning graph quantized tokenizers. In *The Thirteenth*  
998 *International Conference on Learning Representations*, 2025. URL [https://openreview.](https://openreview.net/forum?id=oYSsbY3G4o)  
999 [net/forum?id=oYSsbY3G4o](https://openreview.net/forum?id=oYSsbY3G4o).
- 1000  
1001 Wenjie Wang, Yiyang Xu, Fuli Feng, Xinyu Lin, Xiangnan He, and Tat-Seng Chua. Diffusion rec-  
1002 commender model. In *Proceedings of the 46th International ACM SIGIR Conference on Research*  
1003 *and Development in Information Retrieval*, pp. 832–841, 2023.
- 1004  
1005 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, and  
1006 Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models. *Advances*  
*in neural information processing systems*, 35:24824–24837, 2022.
- 1007  
1008 Jianyu Wei, Shijie Cao, Ting Cao, Lingxiao Ma, Lei Wang, Yanyong Zhang, and Mao Yang. T-  
1009 mac: Cpu renaissance via table lookup for low-bit llm deployment on edge. In *Proceedings of the*  
1010 *Twentieth European Conference on Computer Systems (EuroSys)*, pp. 278–292, 2025.
- 1011  
1012 Yibo Wen, Chenwei Xu, Jerry Yao-Chieh Hu, and Han Liu. Alignab: Pareto-optimal energy align-  
1013 ment for designing nature-like antibodies. *arXiv preprint arXiv:2412.20984*, 2024.
- 1014  
1015 Felix Wu, Amauri Souza, Tianyi Zhang, Christopher Fifty, Tao Yu, and Kilian Weinberger. Sim-  
1016 plifying graph convolutional networks. In *International conference on machine learning*, pp.  
6861–6871. Pmlr, 2019.
- 1017  
1018 Minkai Xu, Alexander S Powers, Ron O Dror, Stefano Ermon, and Jure Leskovec. Geometric latent  
1019 diffusion models for 3d molecule generation. In *International Conference on Machine Learning*,  
pp. 38592–38610. PMLR, 2023.
- 1020  
1021 Weizhi Xu, Junfei Wu, Qiang Liu, Shu Wu, and Liang Wang. Evidence-aware fake news detection  
1022 with graph neural networks. In *Proceedings of the ACM web conference 2022*, pp. 2501–2510,  
1023 2022.
- 1024  
1025 Yuzhuang Xu, Xu Han, Zonghan Yang, Shuo Wang, Qingfu Zhu, Zhiyuan Liu, Weidong Liu, and  
Wanxiang Che. Onebit: Towards extremely low-bit large language models. *Advances in Neural*  
*Information Processing Systems (NeurIPS)*, 37:66357–66382, 2024a.

- 1026 Zhuoyan Xu, Zhenmei Shi, Junyi Wei, Fangzhou Mu, Yin Li, and Yingyu Liang. Towards few-shot  
1027 adaptation of foundation models via multitask finetuning. In *The Twelfth International Conference*  
1028 *on Learning Representations*, 2024b.
- 1029
- 1030 Jiyan Yang, Yin-Lam Chow, Christopher Ré, and Michael W Mahoney. Weighted sgd for  $ell_p$   
1031 regression with randomized preconditioning. *Journal of Machine Learning Research*, 18(211):  
1032 1–43, 2018.
- 1033 Zhengyi Yang, Jiancan Wu, Zhicai Wang, Xiang Wang, Yancheng Yuan, and Xiangnan He. Generate  
1034 what you prefer: Reshaping sequential recommendation via guided diffusion. *Advances in Neural*  
1035 *Information Processing Systems*, 36:24247–24261, 2023.
- 1036
- 1037 Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L. Griffiths, Yuan Cao, and Karthik R  
1038 Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. In *Thirty-*  
1039 *seventh Conference on Neural Information Processing Systems*, 2023.
- 1040 Lu Yi and Zhewei Wei. Scalable and certifiable graph unlearning: Overcoming the approximation  
1041 error barrier. In *The Thirteenth International Conference on Learning Representations*, 2025.  
1042 URL <https://openreview.net/forum?id=pPyJyeLriR>.
- 1043
- 1044 Hao Yu, Yang Zhou, Bohua Chen, Zelan Yang, Shen Li, Yong Li, and Jianxin Wu. Treasures in  
1045 discarded weights for llm quantization. In *Proceedings of the AAAI Conference on Artificial*  
1046 *Intelligence*, volume 39, pp. 22218–22226, 2025.
- 1047 Shen Yuan and Hongteng Xu. Towards better multi-head attention via channel-wise sample permu-  
1048 tation. *arXiv preprint arXiv:2410.10914*, 2024.
- 1049
- 1050 Ofir Zafrir, Guy Boudoukh, Peter Izsak, and Moshe Wasserblat. Q8bert: Quantized 8bit bert. In  
1051 *2019 Fifth Workshop on Energy Efficient Machine Learning and Cognitive Computing-NeurIPS*  
1052 *Edition (EMC2-NIPS)*, pp. 36–39. IEEE, 2019.
- 1053 Manzil Zaheer, Guru Guruganesh, Kumar Avinava Dubey, Joshua Ainslie, Chris Alberti, Santiago  
1054 Ontanon, Philip Pham, Anirudh Ravula, Qifan Wang, Li Yang, et al. Big bird: Transformers for  
1055 longer sequences. *Advances in neural information processing systems*, 33:17283–17297, 2020.
- 1056
- 1057 Amir Zandieh, Insu Han, Majid Daliri, and Amin Karbasi. Kdeformer: Accelerating transformers  
1058 via kernel density estimation. In *ICML*. arXiv preprint arXiv:2302.02451, 2023.
- 1059
- 1060 Chao Zeng, Songwei Liu, Yusheng Xie, Hong Liu, Xiaojian Wang, Miao Wei, Shu Yang, Fangmin  
1061 Chen, and Xing Mei. Abq-llm: Arbitrary-bit quantized inference acceleration for large language  
1062 models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pp. 22299–  
1063 22307, 2025.
- 1064 Yuansong Zeng, Zhuoyi Wei, Zixiang Pan, Yutong Lu, and Yuedong Yang. A robust and scalable  
1065 graph neural network for accurate single-cell classification. *Briefings in Bioinformatics*, 23:2:  
1066 bbab570, 2022.
- 1067
- 1068 Yuchen Zeng and Kangwook Lee. The expressive power of low-rank adaptation. In *The Twelfth*  
1069 *International Conference on Learning Representations*, 2024.
- 1070 Jiahao Zhang. Graph unlearning with efficient partial retraining. In *Companion Proceedings of the*  
1071 *ACM on Web Conference 2024*, pp. 1218–1221, 2024.
- 1072
- 1073 Michael Zhang, Kush Bhatia, Hermann Kumbong, and Christopher Ré. The hedgehog & the porcu-  
1074 pine: Expressive linear attentions with softmax mimicry. In *ICLR*, 2024.
- 1075
- 1076 Renrui Zhang, Jiaming Han, Aojun Zhou, Xiangfei Hu, Shilin Yan, Pan Lu, Hongsheng Li, Peng  
1077 Gao, and Yu Qiao. Llama-adapter: Efficient fine-tuning of language models with zero-init atten-  
1078 tion. *arXiv preprint arXiv:2303.16199*, 2023a.
- 1079
- 1079 Ruiqi Zhang, Spencer Frei, and Peter L Bartlett. Trained transformers learn linear models in-context.  
*arXiv preprint arXiv:2306.09927*, 2023b.

1080 Yanfu Zhang, Shangqian Gao, Jian Pei, and Heng Huang. Improving social network embedding via  
1081 new second-order continuous graph neural networks. In *Proceedings of the 28th ACM SIGKDD*  
1082 *conference on knowledge discovery and data mining*, pp. 2515–2523, 2022.

1083

1084 Ren-Xin Zhao, Jinjing Shi, and Xuelong Li. Qksan: A quantum kernel self-attention network.  
1085 *TPAMI*, 2024.

1086 Zihao Zhao, Eric Wallace, Shi Feng, Dan Klein, and Sameer Singh. Calibrate before use: Improving  
1087 few-shot performance of language models. In *International Conference on Machine Learning*.  
1088 PMLR, 2021.

1089

1090 Huaixiu Steven Zheng, Swaroop Mishra, Xinyun Chen, Heng-Tze Cheng, Ed H Chi, Quoc V Le, and  
1091 Denny Zhou. Step-back prompting enables reasoning via abstraction in large language models.  
1092 In *The Twelfth International Conference on Learning Representations*, 2024.

1093 Chunting Zhou, Pengfei Liu, Puxin Xu, Srini Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat,  
1094 Ping Yu, LILI YU, Susan Zhang, Gargi Ghosh, Mike Lewis, Luke Zettlemoyer, and Omer Levy.  
1095 LIMA: Less is more for alignment. In *Thirty-seventh Conference on Neural Information Process-*  
1096 *ing Systems*, 2023.

1097

1098

1099

1100

1101

1102

1103

1104

1105

1106

1107

1108

1109

1110

1111

1112

1113

1114

1115

1116

1117

1118

1119

1120

1121

1122

1123

1124

1125

1126

1127

1128

1129

1130

1131

1132

1133

# Appendix

**Roadmap.** In Section A, we provide several simple definitions. In Section B, we show more perturbation properties for the softmax matrix. In Section C, we provide the proofs related to network layers with multiple attention heads. In Section D, we prove our main Theorem. In Section F, we provide the broader impact of this work. In Section G, we discuss the LLM usage.

## A PRELIMINARIES

Note that  $\text{softm}(X)$  function defined above is ignoring the effect of the weights  $W_v$ . Here we incorporate them in another function  $\text{softmv}(X)$  (which is also usually called self-attention).

**Definition A.1.** Let  $W_q, W_k$  be weights being used in  $\text{softm}$ . Let  $W_v$  denote the extra weights that will be used in  $\text{softmv}$ . We define  $\text{softmv}(X)$  as follows

$$\text{softmv}(X) := \text{softm}(X)XW_v.$$

*In particular,*

$$\text{softm}(X; W_Q, W_K) := D^{-1} \exp(XW_QW_K^\top X^\top)$$

where  $D$  is diagonal matrix  $D = \text{diag}(\exp(XW_QW_K^\top X^\top)\mathbf{1}_n)$ .

*For notational convenience, we will omit  $W_Q, W_K$  when they are clear from context and simply write  $\text{softm}(X)$ .*

Next we define a very useful parameter  $\epsilon_\ell$  which captures the Lipschitz and layer norm property of every layer.

**Definition A.2.** Let all layers' weights are bounded, i.e.,  $\|W_q\|_\infty, \|W_k\|_\infty, \|W_v\|_\infty \leq \eta$ . Let  $X_0$  denote the first layer input of entire neural network and it is bounded  $\|X_0\|_\infty \leq \phi_0$ . Let  $H$  denote the number of heads. For each layer  $\ell \in [L]$ , we define a parameter  $\epsilon_\ell := 2\eta\phi_0(1 + dH\eta)^\ell$ . We will select parameters to enforce that  $\phi_0 \leq \epsilon_\ell$  in order to simplify some calculations below.

**Definition A.3.** We define function  $g(\epsilon) := 4\epsilon$ .

## B PERTURBATION PROPERTY OF SOFTMAX MATRIX

**Lemma B.1.** If the following conditions hold: Let  $a, b \in \mathbb{R}^n$ . Let  $\|b\|_\infty \leq \epsilon$ . Then we can show

- $|\alpha(a+b)^{-1} - \alpha(a)^{-1}| \leq (e^\epsilon - 1)\alpha(a)^{-1}$
- $|\alpha(a+b)^{-1} - \alpha(a)^{-1}| \leq (e^\epsilon - 1)\alpha(a+b)^{-1}$

*Proof.* We can show that

$$\begin{aligned} |\alpha(a+b)^{-1} - \alpha(a)^{-1}| &= \alpha(a+b)^{-1}\alpha(a)^{-1}|\alpha(a+b) - \alpha(a)| \\ &\leq \alpha(a+b)^{-1}\alpha(a)^{-1} \cdot (e^\epsilon - 1)\alpha(a+b) \\ &= (e^\epsilon - 1)\alpha(a)^{-1} \end{aligned}$$

where the first step follows from simple algebra, the second step follows from Lemma 4.2.

Similarly, we can also show  $|\alpha(a+b)^{-1} - \alpha(a)^{-1}| \leq (e^\epsilon - 1)\alpha(a+b)^{-1}$ .  $\square$

**Lemma B.2.** Let softmax function be defined as Definition 3.1. Let  $a, b \in \mathbb{R}^n$ . If  $|b_i| \leq \epsilon$  for all  $i \in [n]$ , then, we can show that

$$\|\text{softmax}(a+b) - \text{softmax}(a)\|_\infty \leq 2(e^\epsilon - 1)$$

*Proof.* For each  $i \in [n]$ , we can show

$$|\alpha(a+b)^{-1} \exp((a+b)_i) - \alpha(a)^{-1} \exp(a_i)|$$

$$\begin{aligned}
1188 &= |\alpha(a+b)^{-1} \exp((a+b)_i) - \alpha(a+b)^{-1} \exp(a_i) + \alpha(a+b)^{-1} \exp(a_i) - \alpha(a)^{-1} \exp(a_i)| \\
1189 &\leq \alpha(a+b)^{-1} |\exp(b_i + a_i) - \exp(a_i)| + \exp(a_i) \cdot |\alpha(a+b)^{-1} - \alpha(a)^{-1}| \\
1190 &:= A_1 + A_2
\end{aligned}$$

1192 where the second step follows from the triangle inequality.

1193 We can upper bound  $A_1$  as

$$\begin{aligned}
1194 &A_1 = \alpha(a+b)^{-1} \cdot |\exp(a_i + b_i) - \exp(a_i)| \\
1195 &\leq \alpha(a+b)^{-1} \cdot (e^\epsilon - 1) \exp(a_i + b_i) \\
1196 &\leq (e^\epsilon - 1)
\end{aligned}$$

1198 where the second step follows from Lemma 4.2, the third step follows from  $\alpha(x)^{-1} \exp(x_i) \in (0, 1)$  for any  $x$  and  $i$ .

1201 We can upper bound  $A_2$  as

$$\begin{aligned}
1202 &A_2 = \exp(a_i) \cdot |\alpha(a+b)^{-1} - \alpha(a)^{-1}| \\
1203 &\leq \exp(a_i) \cdot (e^\epsilon - 1) \alpha(a)^{-1} \\
1204 &\leq e^\epsilon - 1
\end{aligned}$$

1206 where the second step follows from Lemma B.1, and the third step follows from  $\alpha(x)^{-1} \exp(x_i) \in (0, 1)$  for any  $x$  and  $i$ .

1208 Putting everything together, we can show

$$\begin{aligned}
1209 &\|\text{softmax}(a+b) - \text{softmax}(a)\|_\infty = \max_{i \in [n]} |\alpha(a+b)^{-1} (a+b)_i - \alpha(a)^{-1} a_i| \\
1210 &\leq 2(e^\epsilon - 1).
\end{aligned}$$

1213 Thus, we complete the proof. □

1214 **Lemma B.3.** *If the following conditions hold*

- 1215 • Let  $P = \text{softmax}(A)$  (see Definition 3.1 for function  $\text{softmax}()$ ).
- 1216 • Let  $\tilde{A} = A - E$ .
- 1217 • Let  $\tilde{P} = \text{softmax}(\tilde{A})$ .
- 1218 • Let  $D$  be defined as Definition 3.5, i.e.,  $D_{i,i} := \max_{j,l \in [n]} |E_{i,j} - E_{i,l}|$

1222 Then we can show, for all  $i \in [n], j \in [n]$

$$1224 e^{-D_{i,i}} \tilde{P}_{i,j} \leq P_{i,j} \leq e^{D_{i,i}} \tilde{P}_{i,j}.$$

1226 *Proof.* Let us start by the definition of  $P$ , for each  $i \in [n], j \in [n]$

$$\begin{aligned}
1227 &P_{i,j} = (\text{softmax}(A))_{i,j} \\
1228 &= (\text{softmax}(\tilde{A} + E))_{i,j} \\
1229 &= \frac{\exp(\tilde{A}_{i,j} + E_{i,j})}{\sum_{l=1}^n \exp(\tilde{A}_{i,l} + E_{i,l})} \\
1230 &= \frac{\exp(\tilde{A}_{i,j})}{\sum_{l=1}^n \exp(\tilde{A}_{i,l}) \exp(E_{i,l} - E_{i,j})}
\end{aligned}$$

1236 where the first step follows from definition of  $P$ , the second step follows from definition of  $A$ , the third step follows from definition of softmax (Definition 3.1), and the last step follows from property of exp.

1238 We define  $D_{i,i} := \max_{j,l \in [n]} |E_{i,j} - E_{i,l}|$ . We have that

$$1240 P_{i,j} \in [\tilde{P}_{i,j} \exp(-D_{i,i}), \tilde{P}_{i,j} \exp(D_{i,i})]$$

1241 Thus, we complete the proof. □

1242 **Lemma B.4.** *If the following conditions hold*

- 1243 • Let  $W_q, W_k, W_v$  be the matrix that  $\|W_q\|_\infty, \|W_k\|_\infty, \|W_v\|_\infty \leq \eta$ .
- 1244 • Let  $W = W_q W_k^\top$ .
- 1245 • Let  $E = \beta \text{Res}(X) W \text{Res}(X)^\top$ .
- 1246 • Let  $\beta$  satisfy that  $\beta \leq 1/(\|\text{Res}(X)\|_\infty^2 \eta^2)$ .

1247 Then, we have  $E$  is  $\theta$ -balanced with  $\theta = 1$ .

1248 *Proof.* First, note that  $\|W\|_\infty \leq \|W_q\|_\infty \cdot \|W_k\|_\infty \leq \eta^2$ .

1249 We can show

$$\begin{aligned} \|E\|_\infty &\leq \beta \cdot \|\text{Res}(X)\|_\infty^2 \cdot \|W\|_\infty \\ &\leq \beta \cdot \|\text{Res}(X)\|_\infty^2 \cdot \eta^2 \\ &\leq 1 \end{aligned}$$

1250 Thus, we complete the proof. □

## 1251 C MULTIPLE HEADS

1252 Here, we generalize the proof of Lemma 5.2 to multiple heads. Note that Lemma 5.2 presented a  
1253 simplified proof by ignoring the effects of  $XW_v$ , and thus automatically assuming  $n = d$ . In this  
1254 section we remove that condition and prove the result for general  $n$  and  $d$ . In Section C.1, our goal  
1255 is to prove Lemma C.1, which is the multiple heads version of Lemma 5.2. In Section C.2, we show  
1256 that several required conditions in Lemma C.1 are satisfied.

### 1257 C.1 MULTIPLE HEADS FOR SKIPPING ONE LAYER

1258 In the next Lemma C.1, we will put the effect of  $\text{softmv}$  back. We remark that the major idea of the  
1259 proof remains the same as Lemma 5.2.

1260 **Lemma C.1** (Multiple Heads version of Lemma 5.2). *If the following conditions hold,*

- 1261 • Let  $H$  denote the number of heads.
- 1262 • Note that  $\text{softm}_1$  and  $\text{softm}_2$  are two different instantiations with different  $W_k, W_q, W_v$   
1263 weights.
- 1264 • Let  $X \in \mathbb{R}^{n \times d}$ .
- 1265 •  $A_i = \text{softmv}_{1,i}(X) \in \mathbb{R}^{n \times d}$  for  $i \in [H]$ . (Let  $\text{softmv}$  be defined as Definition A.1)
- 1266 •  $B = X + \sum_{i=1}^H A_i \in \mathbb{R}^{n \times d}$ .
- 1267 •  $\|\text{Res}(A_i)\|_\infty \leq K \cdot \|\text{Res}(X)\|_\infty \leq \epsilon$  for all  $i \in [H]$ . (We remark that this condition will  
1268 hold due to Lemma 5.1; here  $K$  is as defined in Lemma 5.1)
- 1269 • Let  $g(\epsilon) := 4\epsilon$  (see Definition A.3).
- 1270 • Let  $W_v$  satisfy that  $\|(B - X)W_v\|_\infty \leq \epsilon$  and  $n\|XW_v\|_\infty \leq \epsilon \leq 1$  (These conditions will  
1271 be verified by Lemma C.2).
- 1272 • Let  $\epsilon_0 = 3g(3H\epsilon)$ .

1273 Then we can show

- 1274 • **Part 1.**  $\|\text{softm}_2(B) - \text{softm}_2(X)\|_\infty \leq \epsilon_0$
- 1275 • **Part 2.**  $\|\text{softmv}_2(B) - \text{softmv}_2(X)\|_\infty \leq \epsilon_0$

1296 *Proof. Proof of Part 1.*

1297 Let  $R_X = \text{Res}(X)$  so that  $X = R_X + \mathbf{1}_n y_X^\top$  for some vector  $y_X \in \mathbb{R}^d$ .

1299 Using Lemma D.6, we can show that

$$1300 \quad \|\text{softm}_2(X) - \text{softm}_2(R_X)\|_\infty \leq g(\epsilon). \quad (13)$$

1302 To notataionally help in our proof, we define the prefix sums of matrices  $A_0, A_1, \dots, A_i \in \mathbb{R}^{n \times d}$  as

$$1305 \quad A_{[i]} := \sum_{j=0}^i A_j$$

1307 where  $A_0$  is an artificial matrix that has 0 everywhere.

1309 For each  $i \in [H]$ , let  $R_{A_{[i]}} = \text{Res}(A_{[i]})$  so that  $A_{[i]} = R_{A_{[i]}} + \mathbf{1}_n y_{A_{[i]}}^\top$  for some vector  $y_{A_{[i]}} \in \mathbb{R}^d$ .

1311 Using Lemma D.6, we can show that

$$1312 \quad \|\text{softm}_2(A_{[i]}) - \text{softm}_2(R_{A_{[i]}})\|_\infty \leq g(\epsilon).$$

1314 Let  $R_B = \text{Res}(B)$  so that  $B = R_B + \mathbf{1}_n y_B^\top$  for some vector  $y_B \in \mathbb{R}^d$ . Using Lemma D.6, we can show that

$$1316 \quad \|\text{softm}_2(B) - \text{softm}_2(R_B)\|_\infty \leq g(\epsilon) \quad (14)$$

1318 Since  $\|R_{A_{[i]}} - R_{A_{[i-1]}}\|_\infty \leq \epsilon$  for all  $i \in [H]$ , then using Lemma D.6, we have: for each  $i \in [H]$

$$1320 \quad \|\text{softm}_2(R_X + R_{A_{[i]}}) - \text{softm}_2(R_X + R_{A_{[i-1]}})\|_\infty \leq g(\epsilon) \quad (15)$$

1322 We can show  $R_X = \text{Res}(X) = \text{Res}(B - A_{[H]}) = \text{Res}(B - R_{A_{[H]}})$ .

1324 Then, we know

$$1325 \quad \begin{aligned} \|R_X - R_B\|_\infty &= \|\text{Res}(B - R_{A_{[H]}}) - \text{Res}(B)\|_\infty \\ &\leq 2\|R_{A_{[H]}}\|_\infty \end{aligned} \quad (16)$$

1328 where the first step follows from  $R_X = \text{Res}(B - R_{A_{[H]}})$  and  $R_B = \text{Res}(B)$ , the second step follows from Lemma 4.1.

1330 Recall  $B = X + A$ , and  $\|R_A\|_\infty \leq \epsilon$  then we know

$$1332 \quad \begin{aligned} \|R_X + R_A - R_B\|_\infty &\leq \|R_X - R_B\|_\infty + \|R_{A_{[H]}}\|_\infty \\ &\leq 3\|R_{A_{[H]}}\|_\infty \\ &\leq 3H\epsilon, \end{aligned}$$

1336 where the first step follows from triangle inequality, the second step follows from  $\|R_X - R_B\|_\infty \leq 2\|R_{A_{[H]}}\|_\infty$ , and the last step follows from  $\|R_{A_{[H]}}\|_\infty \leq H\epsilon$ .

1338 Since  $\|R_X + R_{A,H} - R_B\|_\infty \leq 3H\epsilon$ , then using Lemma D.6, we have

$$1340 \quad \|\text{softm}_2(R_X + R_{A,H}) - \text{softm}_2(R_B)\|_\infty \leq g(3H\epsilon) \quad (17)$$

1342 Then, we can show

$$1343 \quad \begin{aligned} &\|\text{softm}_2(B) - \text{softm}_2(X)\|_\infty \\ 1344 &\leq \|\text{softm}_2(B) - \text{softm}_2(R_X)\|_\infty + g(\epsilon) \\ 1345 &\leq \|\text{softm}_2(R_B) - \text{softm}_2(R_X)\|_\infty + 2g(\epsilon) \\ 1346 &\leq \|\text{softm}_2(R_B) - \text{softm}_2(R_X + R_{A,H})\|_\infty \\ 1347 &\quad + \sum_{i=1}^{H-1} \|\text{softm}_2(R_X + R_{A,i}) - \text{softm}_2(R_X + R_{A,i-1})\|_\infty + 2g(\epsilon) \end{aligned}$$

$$\begin{aligned}
1350 & \leq g(3H\epsilon) + (H + 2) \cdot g(\epsilon) & (18) \\
1351 & \leq 2g(3H\epsilon) \\
1352 &
\end{aligned}$$

1353 where the first step follows from [triangle inequality and Eq. \(13\)](#), the second step follows from  
1354 [triangle inequality and Eq. \(14\)](#), the third step follows from triangle inequality, the fourth step follows  
1355 from Eq. (15) and Eq. (17), and the last step follows from property of function  $g$ .

### 1356 **Proof of Part 2.**

1357 We can show that

$$\begin{aligned}
1359 & \|\text{softmv}_2(B) - \text{softmv}_2(X)\|_\infty \\
1360 & = \|\text{softm}_2(B)BW_v - \text{softm}_2(X)XW_v\|_\infty \\
1361 & \leq \|\text{softm}_2(B)BW_v - \text{softm}_2(B)XW_v\|_\infty + \|\text{softm}_2(B)XW_v - \text{softm}_2(X)XW_v\|_\infty \\
1362 & \leq \|(B - X)W_v\|_\infty + \|\text{softm}_2(B) - \text{softm}_2(X)\|_\infty \cdot n\|XW_v\|_\infty \\
1363 & \leq \epsilon + 2g(3H\epsilon) \cdot n\|XW_v\|_\infty \\
1364 & \leq 3g(3H\epsilon) \\
1365 &
\end{aligned}$$

1366 where the second step follows from triangle inequality, the third step follows from Fact 3.3 and  
1367 Fact E.2, and the fourth step follows from Eq. (18), where the last step follows from property of  
1368 function  $g$  and assumption in Lemma statement.  $\square$

## 1370 C.2 CONDITIONS IN LEMMA C.1 ARE SATISFIED

1371 Here we will show that the three conditions in Lemma C.1 will be satisfied for each layer  $\ell$ .

- 1372 •  $\|\text{Res}(A_i)\|_\infty \leq K \cdot \|\text{Res}(X)\|_\infty \leq \epsilon$  (where  $K := (e^\theta - 1)\|W_v\|_\infty$ , definition of  $K$  recall  
1373 Lemma 5.1). Here  $\theta = 1$  due to Lemma B.4
- 1374 •  $\|(B - X)W_v\|_\infty \leq \epsilon$
- 1375 •  $n\|XW_v\|_\infty \leq \epsilon$

1376 **Lemma C.2.** *If the following conditions hold*

- 1377 •  $\epsilon_\ell := 2\eta\phi_0(1 + H\eta)^\ell nd$ . (see Definition A.2)
- 1378 • Let  $\eta \in (0, 1]$ .
- 1379 • Let  $\epsilon_\ell \in (0, 1)$ .

1380 Then, we can show

- 1381 • Part 1.  $K \cdot \|\text{Res}(X_\ell)\|_\infty \leq \epsilon_\ell$
- 1382 • Part 2.  $\|(B - X_\ell)W_v\|_\infty \leq \epsilon_\ell$
- 1383 • Part 3.  $n\|X_\ell W_v\|_\infty \leq \epsilon_\ell$

1384 *Proof. Proof of Part 1.*

1385 We can show that

$$\begin{aligned}
1386 & K \cdot \|\text{Res}(X_\ell)\|_\infty \leq 2\eta\|\text{Res}(X_\ell)\|_\infty \\
1387 & \leq 2\eta\|X_\ell\|_\infty \\
1388 & \leq 2\eta\phi_0 \cdot (1 + dH\eta)^\ell \\
1389 & = \epsilon_\ell \\
1390 &
\end{aligned}$$

1391 where the first step follows from  $\theta = 1$  and  $\|W_v\|_\infty \leq \eta$ , the third step follows from Lemma D.8,  
1392 and the last step follows from definition  $\epsilon_\ell$

1393 **Proof of Part 2.**

$$1403 \|(B - X_\ell)W_v\|_\infty \leq \eta \cdot d\|B - X_\ell\|_\infty$$

$$\begin{aligned}
1404 & \\
1405 & = \eta \cdot d \sum_{i=1}^H \|\text{softmv}_i(X_\ell)\|_\infty \\
1406 & \\
1407 & \\
1408 & \leq \eta \cdot d \sum_{i=1}^H \|\text{softm}_i(X_\ell)X_\ell W_{v,i}\|_\infty \\
1409 & \\
1410 & \\
1411 & \leq \eta \cdot d \sum_{i=1}^H \|X_\ell W_{v,i}\|_\infty \\
1412 & \\
1413 & \leq \eta \cdot dH \cdot \|X_\ell W_{v,i}\|_\infty \\
1414 & \leq \eta \cdot dH \cdot \epsilon_\ell \\
1415 & \leq \epsilon_\ell \\
1416 &
\end{aligned}$$

1417 where the first step follows from  $\|W_v\|_\infty \leq \eta$ , the second step follows from definition of  $B$ , the  
1418 third step follows from definition of  $\text{softmv}$ , the fourth step follows from Fact E.2, the sixth step  
1419 follows from part 3, and the last step follows from  $\eta \leq 1/(dH)$ .

### 1420 Proof of Part 3.

1421 We can show

$$\begin{aligned}
1422 & \\
1423 & \|X_\ell W_v\|_\infty \leq \eta d \|X_\ell\|_\infty \\
1424 & \leq \eta d \phi_0 (1 + dH\eta)^\ell \\
1425 & \leq \epsilon_\ell. \\
1426 &
\end{aligned}$$

1427 where the second step follows from Lemma D.8, the third step follows from choice of  $\epsilon_\ell$ . □

## 1430 D MULTIPLE LAYERS

1431  
1432 In Section D.1, we provide the proof of our main theorem. In Section D.2, we provide the Lipschitz  
1433 property of several key functions being used in our proofs. In Section D.3, we prove the Lipschitz  
1434 property for each layer of our Self Attention Network. Finally, in Section D.4, prove the norm of  
1435 each layer in the Self Attention Network is not increasing much.

### 1437 D.1 PROOF OF THEOREM D.1

1438  
1439 **Theorem D.1** (Formal version of Theorem 1.2). *Suppose  $S$  is a SAtt with residuals, with the prop-*  
1440 *erty that for every attention head in every one of its layers, the weight matrices  $W_q, W_k, W_v \in \mathbb{R}^{d \times d}$*   
1441 *all have the bound  $\|W_q\|_\infty, \|W_k\|_\infty, \|W_v\|_\infty \leq \eta$ . Let  $H$  denote the number of heads. Let  $L$  de-*  
1442 *note the number of layers. Assume  $\eta \leq A \cdot \min\{1/(HLd), 1/(\phi_0 nd)\} \leq 1$  for some parameter*  
1443  *$A \leq O(1)$ . Then, there exists a SAtt with residuals  $S'$  with just one layer so that, for any bounded*  
1444  *$X \in \mathbb{R}^{n \times d}$  with  $\|X\|_\infty \leq \phi_0$ , we have  $\|S(X) - S'(X)\|_\infty \leq O(A/L)$ .*

1445 *Proof.* We define

$$1447 \quad X_\ell^{\ell_0} = B_\ell^{\ell_0}$$

1448 Then, we define

$$1449 \quad B_\ell^{\ell_0} = \begin{cases} X_{\ell-1}^{\ell_0} + \sum_{i=1}^H A_{\ell-1,i}^{\ell_0}, & \text{if } \ell \leq \ell_0; \\ \sum_{i=1}^H A_{\ell-1,i}^{\ell_0}, & \text{otherwise.} \end{cases}$$

1453 Let  $\text{softmv}()$  function be defined as Definition A.1. We define

$$1454 \quad A_{\ell-1,i}^{\ell_0} = \text{softmv}_{\ell-1,i}^{\ell_0}(X_{\ell-1}^{\ell_0})$$

1455  
1456 Note that the notation  $B_L^0$  means we have residual in every layer, whereas the notation  $B_L^{\ell_0}$  means  
1457 we don't have a residual connection from layer  $\ell_0$  to layer  $L$ .

Let  $\epsilon_\ell$  be defined as Definition A.2. Let  $\delta := \max_{\ell \in [L]} 3g(3H\epsilon_\ell)$ . Using Lemma C.1, we can show for all  $\ell \in [L]$ ,

$$\|\text{softmv}_\ell(B_\ell^{\ell-1}) - \text{softmv}_\ell(B_\ell^\ell)\|_\infty \leq \delta$$

Let  $C := \max_{\ell \in [L]} 3d\eta(\epsilon_\ell^2 + 1)$ .

Then we can show

$$\begin{aligned} & \|\text{softmv}_2(B_2^0) - \text{softmv}_2(B_2^2)\|_\infty \\ & \leq \|\text{softmv}_2(B_2^0) - \text{softmv}_2(B_2^1)\|_\infty + \|\text{softmv}_2(B_2^1) - \text{softmv}_2(B_2^2)\|_\infty \\ & \leq C \cdot \|\text{softmv}_1(B_1^0) - \text{softmv}_1(B_1^1)\|_\infty + \|\text{softmv}_2(B_2^1) - \text{softmv}_2(B_2^2)\|_\infty \\ & \leq (C + 1)\delta \end{aligned}$$

where the first step follows from triangle inequality, the second step follows from the fact that one layer of the network is  $C$ -Lipschitz (see Lemma D.6), and the last step follows from merging the errors.

For three layers, we have

$$\begin{aligned} & \|\text{softmv}_3(B_3^0) - \text{softmv}_3(B_3^3)\|_\infty \\ & \leq \|\text{softmv}_3(B_3^0) - \text{softmv}_3(B_3^1)\|_\infty \\ & \quad + \|\text{softmv}_3(B_3^1) - \text{softmv}_3(B_3^2)\|_\infty \\ & \quad + \|\text{softmv}_3(B_3^2) - \text{softmv}_3(B_3^3)\|_\infty \\ & \leq C^2 \cdot \|\text{softmv}_1(B_1^0) - \text{softmv}_1(B_1^1)\|_\infty \\ & \quad + C \cdot \|\text{softmv}_2(B_2^1) - \text{softmv}_2(B_2^2)\|_\infty \\ & \quad + \|\text{softmv}_3(B_3^2) - \text{softmv}_3(B_3^3)\|_\infty \\ & \leq C^2\delta + C\delta + \delta \\ & = (C^2 + C + 1)\delta \end{aligned}$$

where the first step follows from triangle inequality, the second step follows from one layer of network is  $C$ -Lipshitz (see Lemma D.6), and the forth step follows from Lemma C.1, and the last step follows from merging the errors.

Therefore for  $L$  layers we have

$$\|\text{softmv}_L(B_L^0) - \text{softmv}_L(B_L^L)\|_\infty \leq (C^L + \dots + C + 1)\delta$$

Thus we complete the proof.

Now, we are ready to analyze the final bound, recall that above we have  $\epsilon_\ell = 2\phi_0\eta nd(1 + dH\eta)^\ell \in (0, 1)$ ,  $\delta = \max_\ell 3g(3H\epsilon_\ell)$ , and  $C = \max_\ell 3d\eta(\epsilon_\ell^2 + 1)$ .

Recall that  $\eta \leq A \cdot \min\{1/(HLd), 1/(\phi_0 nd)\} \leq 1$  for some parameter  $A \leq O(1)$ . Then we can show  $\epsilon_\ell = 2\phi_0\eta nd(1 + dH\eta)^\ell = 2\phi_0\eta nd e^{\ell dH\eta} \leq 2\phi_0\eta nd e^A < 1$ . Next, we can show that compute  $C \leq 3d\eta(\epsilon_\ell^2 + 1) \leq 3d\eta \cdot 2 < 0.5$ . Thus  $C^L + \dots + C + 1 \leq 2$ .

Note that  $3H\epsilon_\ell \in (0, 0.5)$   $\delta \leq 20H\epsilon_\ell \leq AH \cdot \min\{1/(HLd), 1/(\phi_0 nd)\}$

The final is  $2\delta \leq 2AH \cdot \min\{1/(HLd), 1/(\phi_0 nd)\} \leq 2A/L$ .

□

**Remark D.2.** We remark that our proof can be straightforwardly generalized to the situation where the Self Attention Network also has MLP layers, similar to Section 3.2 in (Dong et al., 2021), by defining  $X_\ell^{\ell_0} = f(B_\ell^{\ell_0})$  where  $f$  is the MLP layer. Note that the Lipschitz property of  $f$  will appear correspondingly in the final bound.

**Remark D.3.** Note that in our Theorem D.1, we assume that  $\eta \leq O(\min\{1/(HLd), 1/(\phi_0 nd)\})$ . Meanwhile, the prior work (Dong et al., 2021, Corollary 2.3) similarly requires  $\beta \leq O(\sqrt{d}/(H\phi_0))$ .

1512 Recall that  $\beta$  is in terms of the  $\ell_1$  norm, and so may be up to a factor of  $d^2$  larger than our  $\eta$ ; their  
 1513 factor of  $\sqrt{d}$  only modestly helps with this.

1514  
 1515 (Their result statement is in terms of  $\|\text{Res}(X)\|_\infty$  rather than  $\|X\|_\infty$ , but we could state ours in  
 1516 terms of  $\text{Res}(X)$  instead; we simply bound  $\|\text{Res}(X)\|_\infty \leq \|X\|_\infty$  in the proof of our Lemma C.2.)

## 1517 D.2 LIPSCHITZ PROPERTY

1518 We state a simple application of Lemma B.2.

1519 **Corollary D.4.** Let  $a, b \in \mathbb{R}^n$ . Then, we can show that

$$1520 \|\text{softmax}(a + b) - \text{softmax}(a)\|_\infty \leq 2(e^{\|b\|_\infty} - 1)$$

1521 *Proof.* The proof is same as Lemma B.2. □

1522 **Lemma D.5.** Let  $a, b \in \mathbb{R}^n$ . If  $\|b\|_\infty \leq 1$ , then we have

$$1523 \|\text{softmax}(a + b) - \text{softmax}(b)\|_\infty \leq 4\|b\|_\infty$$

1524 *Proof.* Note that for  $x \in (0, 1]$ , we know  $e^x - 1 \leq 2x$ .

1525 Thus, we know

$$1526 \|\text{softmax}(a + b) - \text{softmax}(b)\|_\infty \leq 2(e^{\|b\|_\infty} - 1) \\ 1527 \leq 4\|b\|_\infty$$

1528 where the first step follows from Corollary D.4, the second step follows from  $e^x - 1 \leq 2x$ . □

1529 **Lemma D.6.** If the following conditions hold

- 1530 • Let  $W_q, W_k, W_v$  denote weight matrices.
- 1531 • Let  $W = W_q W_k^\top$ .
- 1532 • Let  $Y$  satisfy that  $\|Y - X\|_\infty \leq 2\|X\|_\infty$
- 1533 • Let  $d\|X\|_\infty \eta \leq 1/4$ . (This condition is verified in Lemma D.8)
- 1534 •  $K_1 := 12d\|X\|_\infty \|W\|_\infty$ .
- 1535 •  $K_2 := K_1 nd\|X\|_\infty \|W_v\|_\infty + d\|W_v\|_\infty$

1536 Then, we can show

- 1537 • **Part 1.**

$$1538 \|\text{softm}(X) - \text{softm}(Y)\|_\infty \leq K_1 \cdot \|X - Y\|_\infty$$

1539 Further,  $K_1 = 4$ .

- 1540 • **Part 2.**

$$1541 \|\text{softmv}(X) - \text{softmv}(Y)\|_\infty \leq K_2 \cdot \|X - Y\|_\infty$$

1542 *Proof.* Before going to prove each parts, we will first show

$$1543 \begin{aligned} \|XWX^\top - YWY^\top\|_\infty &\leq \|XWX^\top - XWY^\top\|_\infty + \|XWY^\top - YWY^\top\|_\infty \\ &\leq \|XW\|_\infty \cdot \|X - Y\|_\infty + \|WY^\top\|_\infty \cdot \|X - Y\|_\infty \\ &\leq (\|WX\|_\infty + \|WY\|_\infty) \cdot \|X - Y\|_\infty \\ &\leq 3d\|W\|_\infty \|X\|_\infty \|X - Y\|_\infty \\ &\leq 3d\|W\|_\infty 2\|X\|_\infty^2 \\ &\leq 6d^2 \eta^2 \|X\|_\infty^2 \end{aligned}$$

$$\leq 1$$

the first step follows triangle inequality, the second step follows from Fact 3.3, the third step follows from Fact 3.3, the sixth step follows from  $\|W\|_\infty \leq d\eta^2$ , the last step follows from assumption in the Lemma statement.

**Proof of Part 1.** We can show

$$\begin{aligned} \|\text{softm}(X) - \text{softm}(Y)\|_\infty &= \|\text{softmax}(XWX^\top) - \text{softmax}(YWY^\top)\|_\infty \\ &\leq 4\|XWX^\top - YWY^\top\|_\infty \\ &\leq 12d\|W\|_\infty\|X\|_\infty\|X - Y\|_\infty \end{aligned}$$

where the second step follows from Lemma D.5 with  $\|XWX^\top - YWY^\top\|_\infty \leq 1$ .

**Proof of Part 2.** We can show

$$\begin{aligned} &\|\text{softmv}(X) - \text{softmv}(Y)\|_\infty \\ &= \|\text{softm}(X)XW_v - \text{softm}(Y)YW_v\|_\infty \\ &\leq \|\text{softm}(X)XW_v - \text{softm}(Y)XW_v\|_\infty + \|\text{softm}(Y)XW_v - \text{softm}(Y)YW_v\|_\infty \\ &\leq \|\text{softm}(X) - \text{softm}(Y)\|_\infty \cdot \|XW_v\|_\infty \cdot n + \|(X - Y)W_v\|_\infty \\ &\leq K_1 n \|XW_v\|_\infty \|X - Y\|_\infty + d \cdot \|W_v\|_\infty \|X - Y\|_\infty \\ &\leq K_1 nd \|X\|_\infty \|W_v\|_\infty \|X - Y\|_\infty + d \cdot \|W_v\|_\infty \|X - Y\|_\infty \end{aligned}$$

where the first step follows from definition, the second step follows from triangle inequality, the third step follows from Fact 3.3 and Fact E.2, and the fourth step follows from Part 1 and Fact 3.3, and the last step follows from Fact 3.3.

□

1592

### D.3 INSTANTIATING AN INSTANCE FOR EACH LAYER LIPSCHITIZ PROPERTY

1594

**Lemma D.7.** *If the following conditions hold*

- Let  $X_\ell$  denote  $\ell$ -th layer output
- Let  $\|W_q\|_\infty, \|W_k\|_\infty, \|W_v\|_\infty \leq \eta$
- Let  $Y$  satisfy that  $\|Y - X_\ell\|_\infty \leq 2\|X_\ell\|_\infty$
- $\epsilon_\ell := 2\eta\phi_0(1 + dH\eta)^\ell nd$ .

Then, we can show

$$\|\text{softmv}(X_\ell) - \text{softmv}(Y)\|_\infty \leq 3d\eta(\epsilon_\ell^2 + 1)$$

*Proof.* We can show

$$\|\text{softmv}(X_\ell) - \text{softmv}(Y)\|_\infty \leq K_2 \cdot \|X_\ell - Y\|_\infty$$

We just need to upper bound  $K_2$

$$\begin{aligned} K_2 &= K_1 nd \|X_\ell\|_\infty \|W_v\|_\infty + d \|W_v\|_\infty \\ &\leq 12nd^2 \|X_\ell\|_\infty^2 \|W\|_\infty \|W_v\|_\infty + d \|W_v\|_\infty \\ &\leq 12nd^2 \|X_\ell\|_\infty^2 \eta^3 + d\eta \\ &\leq 12nd^2 (\phi_0 \cdot (1 + dH\eta))^{2\ell} \eta^3 + d\eta \\ &= 3d\eta(\epsilon_\ell^2 + 1) \end{aligned}$$

where the first step follows from the definition of  $K_2$ , the fourth step follows from Lemma D.8, and the fifth step follows from the definition of  $\epsilon_\ell$ . □

#### D.4 EACH LAYER NORM IS NOT INCREASING MUCH

**Lemma D.8.** *If the following conditions hold*

- *Let  $X_0$  denote the input of first layer of neural network, and satisfy  $\|X_0\|_\infty \leq \phi_0$*
- *For  $\ell \in [L]$ , we use  $X_\ell$  to denote the  $\ell$ -th layer output*
- *Let  $\|W_v\|_\infty \leq \eta$*

*Then, we can show*

- *Part 1. For any  $\ell$ ,  $\|X_{\ell+1}\|_\infty \leq \|X_\ell\|_\infty \cdot (1 + dH\eta)$*
- *Part 2. For any  $\ell$ ,  $\|X_\ell\|_\infty \leq \phi_0 \cdot (1 + dH\eta)^\ell$*
- *Part 3. For any  $\ell$ ,  $\|X_\ell\|_\infty d\eta \leq 1/4$*

**Proof. Proof of Part 1.**

For any  $\ell$ , we have

$$\begin{aligned}
 \|X_{\ell+1}\|_\infty &= \|X_\ell + \sum_{i=1}^H \text{softmv}_i(X_\ell)\|_\infty \\
 &\leq \|X_\ell\|_\infty + H \cdot \|\text{softmv}_i(X_\ell)\|_\infty \\
 &= \|X_\ell\|_\infty + H \cdot \|\text{softm}_i(X_\ell)X_\ell W_{v,i}\|_\infty \\
 &\leq \|X_\ell\|_\infty + H \cdot \|X_\ell W_{v,i}\|_\infty \\
 &\leq \|X_\ell\|_\infty + H \cdot d \cdot \|X_\ell\|_\infty \|W_{v,i}\|_\infty \\
 &\leq \|X_\ell\|_\infty (1 + dH\eta)
 \end{aligned}$$

where the first step follows from definition of  $X_1$ , the second step follows from triangle inequality, the third step follows from definition of  $\text{softmv}$ , the fourth step follows from Fact E.2, and the fifth step follows from Fact 3.3, and last step follows  $\|W_v\|_\infty \leq \eta$ .

**Proof of Part 2.**

We can show

$$\begin{aligned}
 \|X_\ell\|_\infty &\leq \|X_{\ell-1}\|_\infty (1 + dH\eta) \\
 &\leq \dots \\
 &\leq \|X_0\|_\infty (1 + dH\eta)^\ell \\
 &\leq \phi_0 \cdot (1 + dH\eta)^\ell
 \end{aligned}$$

where the first step follows from Part 1, the third step follows from recursively applying Part 1, and the last step follows from  $\|X_0\|_\infty \leq \phi_0$ .

**Part 3.**

We can show

$$\|X_\ell\|_\infty d\eta \leq \eta d\phi_0 (1 + dH\eta)^\ell \leq \epsilon_\ell \leq 1/4$$

where the second step follows from definition of  $\epsilon_\ell$ , and last step follows from guarantee of  $\epsilon_\ell$ .

Therefore, we complete the proof.  $\square$

## E PROOF OF LEMMA 4.1

**Lemma E.1** (Restatement of Lemma 4.1). *Let  $\text{Res}()$  be defined as Definition 3.4. If  $\|A - B\|_\infty \leq \epsilon$ , then  $\|\text{Res}(A) - \text{Res}(B)\|_\infty \leq 2\epsilon$ .*

1674 *Proof.* Let  $\mathbf{1}$  denote a column vector where all the entries are ones.

1675  $\text{Res}(Z)$  acts columnwise: if  $z^{(j)}$  is column  $j$  of  $Z$ , then the  $j$ -th column of  $\text{Res}(Z)$  is

$$1677 \quad \text{res}(z^{(j)}) := z^{(j)} - t(z^{(j)})\mathbf{1}$$

1678 where  $t(z^{(j)})$  minimizes  $\|z^{(j)} - t \cdot \mathbf{1}\|_\infty$ .

1680 So it suffices to prove that for vectors  $x, y \in \mathbb{R}^n$  we have

$$1681 \quad \|\text{res}(x) - \text{res}(y)\|_\infty \leq 2 \cdot \|x - y\|_\infty.$$

1683 For any vector  $z \in \mathbb{R}^n$ , let  $M(z) = \max_i z_i$ ,  $m(z) = \min_i z_i$ .

1684 The scalar  $t$  minimizing  $\|z - t\mathbf{1}\|_\infty = \max_i |z_i - t|$  is

$$1686 \quad t(z) := \frac{M(z) + m(z)}{2}$$

1688 so  $\text{res}(z) = z - t(z)\mathbf{1}$ .

1689 Let  $\delta := \|x - y\|_\infty$ . Then

$$1690 \quad |M(x) - M(y)| \leq \delta, |m(x) - m(y)| \leq \delta$$

1692 so

$$1693 \quad |t(x) - t(y)| = \left| \frac{1}{2}(M(x) + m(x) - M(y) - m(y)) \right|$$

$$1694 \quad \leq \frac{1}{2}(\delta + \delta)$$

$$1695 \quad = \delta.$$

1698 Now

$$1699 \quad \text{res}(x) - \text{res}(y) = (x - y) + (t(x) - t(y))\mathbf{1},$$

1700 so for each  $i$ ,

$$1701 \quad |\text{res}(x)_i - \text{res}(y)_i| \leq |x_i - y_i| + |t(x) - t(y)| \leq \delta + \delta = 2\delta$$

1703 Hence

$$1704 \quad \|\text{res}(x) - \text{res}(y)\|_\infty \leq 2 \cdot \|x - y\|_\infty.$$

1706 Applying this to each column of  $A, B$  shows

$$1707 \quad \|\text{Res}(A) - \text{Res}(B)\|_\infty \leq 2 \cdot \|A - B\|_\infty.$$

1708 Thus, we complete the proof.  $\square$

1709 **Fact E.2.** Given  $A \in \mathbb{R}^{n \times n}$  and  $B, C \in \mathbb{R}^{n \times d}$ , we have  $\|\text{softmax}(A)(B - C)\|_\infty \leq \|B - C\|_\infty$

1712 *Proof.* We just need to prove for the case: one row of  $A$  and one column of  $B, C$ .

$$1713 \quad |(\text{softmax}(a), b - c)| \leq \sum_{i=1}^n \text{softmax}(a)_i |b_i - c_i| \leq \sum_{i=1}^n \text{softmax}(a)_i \|b - c\|_\infty = \|b - c\|_\infty$$

1716 Thus, we complete the proof.  $\square$

## 1718 F BROADER IMPACT

1720 Our results offer new theoretical insights into the expressiveness of attention mechanisms in trans-  
1721 formers. These findings may guide the future design of large language models toward more expres-  
1722 sive architectures. We do not foresee any potential negative societal impacts from this work.

## 1724 G LLM USAGE DISCLOSURE

1726 LLMs were used only to polish language, such as grammar and wording. These models did not  
1727 contribute to idea creation or writing, and the authors take full responsibility for this paper’s content.