ATTENTION-ONLY TRANSFORMERS VIA UNROLLED SUBSPACE DENOISING

Anonymous authors

Paper under double-blind review

ABSTRACT

Despite the great success of transformers in practice, their architectures have been empirically designed, hence lack of mathematical justification and interpretability. Moreover, many empirical studies have indicated that some components of the transformer architectures may be redundant and can be removed or replaced without compromising overall performance. Hence to derive a compact and interpretable transformer architecture, we contend that the goal of representation learning is to compress a set of noisy initial token representations towards a mixture of low-dimensional subspaces. Based on the existing literature, the associated denoising operation naturally takes the form of a multi-subspace self-attention (MSSA). By unrolling such iterative denoising operations as a deep network, we arrive at a highly compact architecture that consists of only an MSSA operator with skip connections at each layer, without MLP. We rigorously prove that each layer of the proposed transformer performs so highly efficient denoising that it improves the signal-to-noise ratio of token representations at a linear rate with respect to the number of layers. Despite its simplicity, extensive experiments on language and vision tasks demonstrate that such a minimalistic attention-only transformer can achieve performance close to conventional transformers, such as GPT-2 and CRATE.

027 028 029

004

006

008 009

010 011

012

013

014

015

016

017

018

019

021

024

025

026

1 INTRODUCTION

031 Over the past years, transformer architectures (Vaswani et al., 2017) have achieved remarkable empirical success across various modern machine learning applications, including large language mod-033 els (LLMs) (Devlin, 2018; Brown et al., 2020a), vision generative models (Chen et al., 2020; Bao 034 et al., 2023; Peebles & Xie, 2023), and reinforcement learning (Chen et al., 2021). In general, transformer architectures are constructed by stacking multiple identical layers that work together to process and learn from data. Each layer is composed of several interacting components arranged in a specific sequence, including self-attention operators, layer normalization, multilayer perceptron 037 (MLP) networks, and skip connections. In practice, transformer architectures, such as BERT (Devlin, 2018) and GPT-4 (Achiam et al., 2023), are highly deep, often with dozens to even hundreds of layers, and are significantly over-parameterized, containing millions or even billions of parame-040 ters. This considerable depth and a large number of parameters endow transformers with impressive 041 learning capabilities, allowing them to model complex patterns and relationships in real-world data. 042

Despite the remarkable success of transformers, their deep and over-parameterized architecture 043 makes them complex "black box", hindering interpretability and the understanding of their inner 044 mechanism. To address this, a common approach involves systematically removing or modifying certain components in transformers to simplify the architecture; see, e.g., Dong et al. (2021); Alcalde 046 et al. (2024); Noci et al. (2024); Geshkovski et al. (2023a); Geva et al. (2020); Guo et al. (2024). For 047 example, Alcalde et al. (2024) studied pure-attention hard-max transformers with skip connections 048 and showed that the output converges to a clustered equilibrium as the number of layers goes to infinity. Noci et al. (2024) analyzed a modified softmax-based attention model with skip connections, demonstrating that the limiting distribution can be described by a stochastic differential equation. 051 These studies indicate that the most basic components of transformers are self-attention layers and skip connections. Although existing studies have provided valuable insights into different com-052 ponents of transformers, few of them elucidate the underlying mechanisms by which transformers process and transform input into output across layers.

054 Moreover, existing empirical studies suggest that some components of transformers are not be es-055 sential and can be removed or modified without compromising performance. For example, He & 056 Hofmann (2024) empirically demonstrated that transformer architecture can be simplified by remov-057 ing components such as skip connections, value matrix, and normalization layers without degrad-058 ing performance. Additionally, Sukhbaatar et al. (2019) investigated the effects of removing MLP blocks from transformers and augmenting the self-attention layers to play a similar role to MLP blocks, showing that performance can be preserved. Similarly, Pires et al. (2023) examined the po-060 tential for reducing the frequency of MLP layers in transformers. Other works also studied other 061 simplifications of transformers, such as linear attentions (Katharopoulos et al., 2020) and shared-062 QK attentions (Kitaev et al., 2020). Based on these discussions, this work focuses on addressing the 063 following question regarding the understanding of the underlying mechanism of transformers and 064 the design of their architectures:

Can we design a minimalistic transformer-like deep architecture consisting of fully interpretable

and provably effective layers that achieves performance close to that of standard transformers?

065 066

067

068

084

069 1.1 RELATED WORKS

Existing studies on self-attention mechanisms. It is widely believed that the power of transform-071 ers primarily stems from their self-attention layers, which enable the model to capture long-range 072 dependencies and contextual relationships between tokens by dynamically weighing token relation-073 ships across the input sequence (Tsai et al., 2019; Vaswani et al., 2017). To explore the mechanism 074 behind self-attention, numerous studies have investigated the performance of pure self-attention networks, often incorporating only one additional component to prevent rank collapse and maintain 075 expressiveness. For example, Dong et al. (2021) showed that in pure-attention transformers without 076 skip connections and MLP layers, token representations collapse exponentially to a rank-1 matrix 077 across layers. They also showed that self-attention networks with skip connections prevent rank collapse. Geshkovski et al. (2023a;b) studied the dynamics of multi-head self-attentions and char-079 acterized clustering behaviors of learned representations. Recently, Wu et al. (2024) showed that pure self-attention networks with LayerNorm can prevent rank collapse. While these studies have 081 advanced the theoretical understanding of self-attention mechanisms in simplified transformer archi-082 tectures, they don't provide any empirical validation on real-world vision or language tasks, offering 083 little insight into the role of self-attention in practice.

Deep network architecture design via unrolled optimization. 085 It is commonly believed that the success of modern deep networks 086 largely stems from their ability to transform the raw data into 087 compact and structured representations, which facilitates down-880 stream tasks (Chan et al., 2022; Chen et al., 2023; Ma et al., 089 2022; Yu et al., 2023a). A principled and interpretable approach 090 to learning such representations with transformers is to construct 091 an architecture that incrementally transforms tokens into these 092 representations via unrolling iterative optimization steps as layers of a deep network (Chan et al., 2022; Monga et al., 2021; Wang et al., 2016; Yu et al., 2023b; Zhang & Ghanem, 2018). 094 Notably, Monga et al. (2021) demonstrate that such unrolled net-095 works are more interpretable, parameter-efficient, and effective 096 compared to generic networks. In this approach, each iteration of an algorithm for learning compact and structured representa-098 tions is represented as one layer of deep networks. For example, Gregor & LeCun (2010) have demonstrated that sparse coding al-100



Figure 1: Each layer of the proposed attention-only transformer architecture.

gorithms, such as ISTA, can be used to construct MLPs. Recently, Chan et al. (2022) constructed 101 a "white-box" network based on an iterative gradient descent scheme to optimize the maximal cod-102 ing rate reduction objective. More recently, Yu et al. (2023a) designed a "white-box" transformer 103 architecture by implementing an approximate alternating minimization to optimize the sparse rate 104 reduction objective. The proposed transformer achieves performance comparable to some popular 105 ones such as ViT (Dosovitskiy et al., 2020), BERT (Devlin, 2018), and DINO (Caron et al., 2021) on vision tasks. Notably, a key component in their design is the multi-head subspace self-attention 106 (MSSA) operator (see Eq. (3)). While they argued that this operator can denoise token representa-107 tions, they only showed that the negative gradient of the compression term of the objective points to the denoising direction, without providing an accurate analysis or guarantee for the denoising efficiency. The MSSA's denoising capabilities remain an open question.

Linear representation & superposition hypotheses. Recent empirical studies on language tasks 111 have raised the "linear representation hypothesis", which posits that token representations can be 112 linearly encoded as one-dimensional feature vectors in the activation space of LLMs (Jiang et al., 113 2024; Park et al., 2023), and "superposition hypothesis", which further hypothesizes that token rep-114 resentations are a sparse linear combination of these feature vectors (Elhage et al., 2022; Yun et al., 115 2021; Arora et al., 2018). Building on these hypotheses, various approaches have been proposed 116 to understand and utilize token representations. For example, Templeton (2024) employed sparse 117 autoencoders to decompose the token representations of Claude 3 Sonnet into more interpretable pieces. Luo et al. (2024) leveraged sparse dictionary learning to explore token representations, de-118 composing them into interpretable components based on a concept dictionary. Recently, Engels 119 et al. (2024) conjectured that token representations in LLMs are the sum of many sparse multi-120 dimensional features. This conjecture is supported by their experiments on GPT-2 and Mistral 7B, 121 where they used sparse autoencoders to identify multi-dimensional features. Notably, all of these 122 empirical studies come to the qualitative conclusion that the token representations lie on a union of 123 (possibly many) low-dimensional subspaces. 124

125 1.2 OUR CONTRIBUTIONS

126 Based on the above discussions, we use a simple yet effective model for the token representations 127 that accurately reflects the behaviors of trained transformers (such as LLMs) based on the previously 128 referenced empirical studies. That is, we model the underlying distribution of token representations as a mixture of low-rank Gaussians corrupted by noise (see Definition 1). Specifically, each token 129 representation lies in a subspace corrupted by the noise from other spaces (see Eq. (1)). To denoise 130 these token representations, we employ the multi-head subspace self-attention (MSSA) operator 131 proposed in (Yu et al., 2023a; Pai et al., 2023) to incrementally update the token representations (see 132 Eq. (3)). Then, our contributions can be summarized as follows: 133

133 134 135

136

137

138

• Attention-only transformer with a minimalistic architecture via unrolled optimization. Based on unrolling the iterative optimization steps Eq. (3), we construct a new transformer with a streamlined architecture, consisting of only MSSA layers with skip connections (see Figure 1).¹ This design simplifies transformer architectures significantly compared to standard decoder-only transformers. More details are illustrated in Figure 3.

- Theoretical guarantees on the denoising performance of the proposed transformer. To quantify the denoising performance, we define a signal-to-noise (SNR) metric (see Eq. (8)) for each block of the token representations. We prove that each layer of the proposed transformer improves the SNR at a linear rate when the initial token representations are sampled from a mixture of low-rank Gaussians (see Theorem 1). This indicates the MSSA operator is highly effective in denoising token representations towards their corresponding subspaces.
- Understanding roles of self-attention and MLP layers. Notably, the proposed transformer is a valuable model for understanding the mechanism of attention since it disentangles the effect of MLP layers. Moreover, comparing the proposed transformer to standard transformers provides insights into the specific role, or empirical benefits, of the MLP layers in different tasks, such as for in-context learning (see experiments in Section 4.1.2).

We have conducted extensive experiments on both language and vision tasks, including causal language modeling, in-context learning, and supervised image classification, to complement our theory and demonstrate the potential of our proposed transformer architecture. These experiments highlight its ability to handle complex real-world applications, thereby confirming the practical value of our streamlined attention-only transformer design.

Notation. Given an integer n, we denote by [n] the set $\{1, \ldots, n\}$. Given a vector a, let ||a|| denote the Euclidean norm of a and diag(a) denote the diagonal matrix with a as its diagonal. Given a matrix A, let ||A|| denote the spectral norm of A, $||A||_F$ denote the Frobenius norm, and a_{ij} denote the (i, j)-th element. For sequences of positive numbers $\{a_n\}$ and $\{b_n\}$, we write $a_n \leq b_n$ or $b_n \geq a_n$ if there exists an absolute constant C > 0 such that $a_n \leq Cb_n$. Given a constant $\tau > 0$, we define $\mathbb{I}(x > \tau) = 1$ if $x > \tau$ and $\mathbb{I}(x > \tau) = 0$ otherwise.

¹For language tasks, we additionally include LayerNorm layers to improve performance.

162 2 TECHNICAL APPROACH AND JUSTIFICATION

164 To begin, we introduce the basic setup of transformers for learning representations from real-world 165 data. Real-world data, such as images, videos, and text, are often modeled as random samples 166 drawn from a high-dimensional probability distribution with low-dimensional intrinsic structures 167 (Pope et al., 2020; Wright & Ma, 2022). Instead of directly inputting data samples into trans-168 formers, a common preprocessing step involves converting each sample into a sequence of vectors, referred to as tokens. Each token represents a localized segment of the data, such as an image patch, 169 a snippet of text, or a frame in a video. Consequently, the input to transformers is typically a se-170 quence of tokens, denoted as $X = [x_1, \dots, x_N] \in \mathbb{R}^{D \times N}$. Then, the goal of transformers is to 171 learn a map $f: \mathbb{R}^{D \times N} \to \mathbb{R}^{d \times N}$ that transforms these tokens into structured and compact token 172 representations that facilitate downstream tasks, such as classification (Dosovitskiy et al., 2020), 173 segmentation (Kirillov et al., 2023), and generation (Saharia et al., 2022), by capturing the underly-174 ing patterns and relationships in the data. For ease of exposition, we denote the token representations 175 as $\boldsymbol{Z} := f(\boldsymbol{X}) \in \mathbb{R}^{d \times N}$.

176 177

178

187 188

194 195

196

200

201

2.1 LEARNING TOKEN REPRESENTATIONS VIA UNROLLED OPTIMIZATION

In this subsection, we introduce how to learn token representations based on the approach of unrolling optimization algorithms (Chan et al., 2022; Gregor & LeCun, 2010; Monga et al., 2021; Sun et al., 2019; Wang et al., 2016; Yu et al., 2023b; Zhang & Ghanem, 2018). This approach involves constructing each layer of a neural network according to a step of an iterative optimization algorithm. That is, the network's architecture is designed to implement a specific optimization algorithm, where each layer corresponds to a single iterative step. By unrolling the algorithm, a "white-box" transformer architecture can be constructed as a multi-layer neural network that incrementally transforms input tokens into structured and compact representations. This process can be described as follows:

$$f: \mathbf{X} \xrightarrow{f^0} \mathbf{Z}^{(0)} \xrightarrow{f^1} \cdots \xrightarrow{f^l} \mathbf{Z}^{(l)} \xrightarrow{f^{l+1}} \cdots \xrightarrow{f^L} \mathbf{Z}^{(L)} = \mathbf{Z},$$

where $f^0 : \mathbb{R}^{D \times N} \to \mathbb{R}^{d \times N}$ is a pre-processing mapping (e.g., positional encoding, token embedding) that transforms input tokens $X \in \mathbb{R}^{D \times N}$ to initial token representations $Z^{(0)} \in \mathbb{R}^{d \times N}$, $f^l : \mathbb{R}^{d \times N} \to \mathbb{R}^{d \times N}$ denotes an incremental operation, and $Z^{(l)}$ denotes the token representations at the *l*-th layer for each $l \in [L]$. Then, a key question is how to design the operator f^l at each layer to learn meaningful token representations efficiently throughout the network in a principled manner.

2.2 DENOISING OPERATOR FOR LEARNING TOKEN REPRESENTATIONS

In this subsection, we introduce a denoising operator for learning token representations incrementally. To clarify the intuition behind this design, we assume that the initial token representations $Z^{(0)}$ are drawn from a mixture of noisy low-rank Gaussian distributions as follows.

Definition 1. Let C_1, \ldots, C_K be a partition of the index set [N] and $U_k \in \mathbb{R}^{d \times p_k}$ denote the orthonormal basis of the k-th cluster for each $K \in [K]$. We say that the token representations $\{z_i^{(0)}\}_{i=1}^N$ are sampled from a mixture of noisy low-rank Gaussian distributions if for each $k \in [K]$,

206

$$\boldsymbol{z}_{i}^{(0)} = \boldsymbol{U}_{k}\boldsymbol{a}_{i} + \sum_{j \neq k}^{K} \boldsymbol{U}_{j}\boldsymbol{e}_{i,j}, \quad \forall i \in C_{k},$$
(1)

where $\mathbf{a}_i \stackrel{i.i.d.}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}_{p_k})$ and $\mathbf{e}_{i,j} \stackrel{i.i.d.}{\sim} \mathcal{N}(\mathbf{0}, \delta^2 \mathbf{I}_{p_j})$ for all $i \in C_k$ and $k \in [K]$, $\{\mathbf{a}_i\}$ and $\{\mathbf{e}_{i,j}\}$ are respectively mutually independent, and $\{\mathbf{a}_i\}$ is independent of $\{\mathbf{e}_{i,j}\}$.

Before proceeding, we make some remarks on this model. First, it provides a probabilistic framework for modeling token representations, assuming that they are sampled from a mixture of multiple low-rank Gaussian distributions with noise. Specifically, if a token representation belongs to the *k*th cluster as shown in Eq. (1), it consists of a signal component $U_k a_i$ and a noise component $\sum_{j \neq k}^{K} U_j e_{i,j}$. Second, this model aligns well with the "linear representation hypothesis" (Jiang et al., 2024; Park et al., 2023) and "superposition hypothesis" (Elhage et al., 2022; Yun et al., 2021; Arora et al., 2018) regarding the structures of token representations in pretrained LLMs. Indeed, the



Figure 2: Layers of transformers f^l gradually denoise token representations towards their corresponding subspaces.

bases of subspaces can be interpreted as semantics features, and each token representation can be approximately expressed as a sparse linear combination of subspace bases when the noise variance δ is sufficiently small. Our goal is to denoise these token representations towards the corresponding subspace; see Figure 2.

Denoising operator for token representations. In this work, we make the simplifying assump-232 tion that the subspaces are orthogonal to each other in Definition 1, i.e., $U_k^T U_j = 0$ for all $k \neq j$. 233 Note this assumption is not so limiting as in high-dimensional spaces, with high-probability low-234 dimensional subspaces are incoherent, i.e., $U_k^T U_j \approx 0$ to each other (Wright & Ma, 2022).² 235

236 Without loss of generality, we rearrange the token representations $Z^{(0)}$ such that the token repre-237 sentations from the same subspace are concatenated together, i.e.,

$$\boldsymbol{Z}^{(0)} = \begin{bmatrix} \boldsymbol{Z}_1^{(0)} & \dots & \boldsymbol{Z}_K^{(0)} \end{bmatrix} = \begin{bmatrix} \boldsymbol{U}_1 \boldsymbol{A}_1 + \sum_{j \neq 1} \boldsymbol{U}_j \boldsymbol{E}_{1,j} & \dots & \boldsymbol{U}_K \boldsymbol{A}_K + \sum_{j \neq K} \boldsymbol{U}_j \boldsymbol{E}_{K,j} \end{bmatrix},$$

where the columns of $Z_k^{(0)}$ denote the token representations from the k-th subspace for each $k \in [K]$, the columns of $A_k \in \mathbb{R}^{p_k \times N_k}$ consists of $\{a_i\}_{i \in C_k}$, and the columns of $E_{k,j} \in \mathbb{R}^{p_j \times N_k}$ consists of $\{e_{i,j}\}_{i \in C_k}$ for each $k \in [K]$ with $N_k = |C_k|$ for each $k \in [K]$. Obviously, projecting 241 242 243 244 token representations onto their corresponding subspace helps separate the signal from the noise 245 components, i.e.,

246 247

225

226 227

228

229

230

231

238 239 240

248

251

 $\boldsymbol{U}_{k}\boldsymbol{U}_{k}^{T}\boldsymbol{Z}_{s}^{(0)} = \begin{cases} \boldsymbol{U}_{k}\boldsymbol{A}_{k}, & \text{if } s = k, \\ \boldsymbol{U}_{k}\boldsymbol{E}_{s\,k}, & \text{if } s \neq k. \end{cases}$ (2)

249 To denoise the token representations from k-th subspace, we can compute the similarity of projected 250 token representations via $(U_k^T Z)^T (U_k^T Z)$ and verify that the similarity between projected token representations from the k-th subspace is high, while the similarity between other pairs of projected 252 token representations is low when $\delta < 1$. Then, we convert it to a distribution of membership with 253 function φ , such as hard-thresholding or soft-max functions, and denoise the token representations 254 towards to the corresponding subspace using this membership. Now, we formalize the considered 255 operator as follows:

256 257

258 259

262 263

269

$$\boldsymbol{Z}^{(l+1)} = \boldsymbol{Z}^{(l)} + \eta \sum_{k=1}^{K} \boldsymbol{U}_{k} \boldsymbol{U}_{k}^{T} \boldsymbol{Z}^{(l)} \varphi \left(\boldsymbol{Z}^{(l)^{T}} \boldsymbol{U}_{k} \boldsymbol{U}_{k}^{T} \boldsymbol{Z}^{(l)} \right), \quad l = 0, 1, \dots, L-1, \quad (3)$$

where $\eta > 0$ is the denoising strength and $\varphi(\cdot) : \mathbb{R}^{d \times N} \to \mathbb{R}^{d \times N}$ is an operator applied to each 260 261 column of an input matrix, i.e.,

$$\varphi(\mathbf{X}) = [\varphi(\mathbf{x}_1) \quad \dots \quad \varphi(\mathbf{x}_N)].$$
 (4)

264 Notably, this operator, referred to as the multi-head subspace self-attention (MSSA), is first proposed 265 by Yu et al. (2023a;b) to approximately optimize the compression term of the sparse rate reduction 266 objective for constructing a transformer architecture. They showed that the negative compression 267 gradient of the objective points from the token representation to the corresponding subspace. However, they do not give any accurate analysis of the denoising efficiency of the MSSA operator (3). 268

²It is not difficult to generalize our results to the more general case, with slightly more sophisticated analysis.



Figure 3: **Details of the attention-only transformer (AoT) architecture.** Each layer consists of the MSSA operator and a skip connection. In addition, LayeNnorm is included only for language tasks. In practice, backpropagation is applied to train the model parameters using training samples.

2.3 TRANSFORMER ARCHITECTURE DESIGN VIA UNROLLED OPTIMIZATION

Now, we formally introduce the proposed transformer architecture. Specifically, by unrolling the iterative optimization steps (3) as layers of a deep network, we construct a transformer architecture in Figure 3. Each layer of the proposed architecture only consists of the MSSA operator and a skip connection. For language tasks, we additionally incorporate LayerNorm before the MSSA operator to improve performance. The complete architecture is built by stacking such layers, along with essential task-specific pre-processing and post-processing steps, such as positional encoding, token embedding, and final task-specific head to adapt to different applications.

Remark 1. Generally speaking, the standard decoder-only transformer architecture is composed of the following key components (Brown et al., 2020b; Radford et al., 2019; Vaswani et al., 2017):
 (1) positional encoding, (2) multi-head QKV self-attention mechanisms, (3) feed-forward MLP networks, (4) layer normalization, and (5) residual connections. In contrast, our proposed transformer architecture adopts a streamlined design by incorporating several key simplications. Specifically, it employs shared-QKV subspace self-attention mechanisms, excludes MLP layers, and reduces the frequency of LayerNorm.

Differences from previous works on attention-only transformers. In the literature, some theo-302 retical works have studied attention-only transformers. For example, Dong et al. (2021); Wu et al. 303 (2024) showed that pure-attention transformers with skip connections or LayerNorm can prevent 304 rank collapse. Additionally, Alcalde et al. (2024) studied the clustering behavior of attention-only 305 hardmax transformers. While these studies contribute significantly to our understanding of the role 306 of self-attention in transformers, they lack empirical validation and practical implications. In con-307 trast to these works, we not only show that each layer of the proposed attention-only transformer 308 can denoise token representations but also conduct experiments on real-world language and vision 309 tasks to demonstrate the potential.

The role of backward propagation. Notably, our approach constructs a transformer architecture in the forward pass by interpreting each layer as a denoising operator, conditioned on the assumption that the subspace bases $\{U_k\}_{k=1}^K$ are known. However, in practice, these subspace matrices, i.e., network parameters, are unknown and need to be learned gradually via iterative optimization too. Hence, the forward denoising operator (3) at the *l*-th layer/iteration becomes

$$\boldsymbol{Z}^{(l+1)} = \boldsymbol{Z}^{(l)} + \eta \sum_{k=1}^{K} \boldsymbol{U}_{k}^{(l)} \boldsymbol{U}_{k}^{(l)^{T}} \boldsymbol{Z}^{(l)} \varphi \left(\boldsymbol{Z}^{(l)^{T}} \boldsymbol{U}_{k}^{(l)} \boldsymbol{U}_{k}^{(l)^{T}} \boldsymbol{Z}^{(l)} \right), \quad l = 0, 1, \dots, L - 1.$$
(5)

317 318 319

316

282

283

284 285

286

301

We should emphasize that the parameters $\{U_k^{(l)}\}$ now depend on the layer index l and can be different across layers. Note that $U_k^{(l)}$ at different layers can represent different intermediate estimates for U_k via certain optimization. In practice, they can be estimated through end-to-end training via backpropagation. This flexibility brings additional capacity for the overall deep architecture, allowing learning denoising bases $\{U_k^{(l)}\}$ at each layer that is locally adaptive to the distribution of $Z^{(l)}$.


Figure 4: **Denosing performance of the attention-only transformer.** Here, we sample initial token representations from a mixture of low-rank Gassuains in Definition 1. Then, we apply (3) to update token representations and report the SNR at each layer.

3 THEORETICAL GUARANTEE FOR THE ATTENTION-ONLY TRANSFORMER

In this section, we rigorously show that each layer of the proposed transformer denoises token representation when the initial token representations are sampled from a mixture of low-rank Gaussians as defined in Definition 1. To quantify the denoising capability, we define the signal-to-noise ratio (SNR) for each block of the token representations at the *l*-th layer as follows:

$$SNR(\boldsymbol{Z}_{k}^{(l)}) := \frac{\|\boldsymbol{U}_{k}\boldsymbol{U}_{k}^{T}\boldsymbol{Z}_{k}^{(l)}\|_{F}}{\|(\boldsymbol{I} - \boldsymbol{U}_{k}\boldsymbol{U}_{k}^{T})\boldsymbol{Z}_{k}^{(l)}\|_{F}}, \quad \forall k \in [K].$$
(6)

To simplify our analysis, we assume that $p = p_1 = \cdots = p_K$, $N_1 = \cdots = N_K = N/K$, and

$$\begin{bmatrix} \boldsymbol{U}_1 & \dots & \boldsymbol{U}_K \end{bmatrix} \in \mathcal{O}^{d \times Kp}. \tag{7}$$

With the above setup, we now characterize the denoising performance of the proposed transformer. **Theorem 1.** Let $Z^{(0)}$ be defined in Definition 1 and $\varphi(\cdot)$ in Eq. (4) be $\varphi(x) = h(\sigma(x))$, where $\sigma : \mathbb{R}^N \to \mathbb{R}^N$ is the soft-max function and $h : \mathbb{R}^N \to \mathbb{R}^N$ is an element-wise thresholding function with $h(x) = \tau \mathbb{I} \{x > \tau\}$ for each $i \in [N]$. Suppose that $p \gtrsim \log N$, $\delta \lesssim \sqrt{\log N} / \sqrt{p}$, and

$$\tau \in \left(\frac{1}{2}, \frac{1}{1+N\exp(-9p/32)}\right].$$

For sufficiently large N, it holds with probability at least $1 - KLN^{-\Omega(1)}$ that for each $l \in [L-1]$,

$$\operatorname{SNR}(\boldsymbol{Z}_{k}^{(l+1)}) = (1 + \eta\tau)\operatorname{SNR}(\boldsymbol{Z}_{k}^{(l)}), \quad \forall k \in [K].$$
(8)

The proof is deferred to Appendix A. Here we comment on significance of this theorem:

- Linear denoising performance of the attention-only transformer. In the theorem, when the initial token representations are sampled from a mixture of low-rank Gaussian distributions with a noise level $O(\sqrt{\log N}/\sqrt{p})$ and $\varphi(\cdot)$ is defined in (4), we show that each layer of the proposed transformer denoises token representations at a linear rate. This indicates the MSSA operator's efficiency in reducing noise across layers. Notably, our theoretical results are well-supported by experimental observations in Figure 4, which further validate the practical denoising capability of the proposed transformer.
- • Difficulties in analyzing the dynamics of the update (3). It is worth noting that the update (3) is highly nonlinear and complicated. Specifically, it is cubic in terms of update variables $Z^{(l)}$ and the operator φ is nonlinear, being composed of soft-max and thresholding functions. These characteristics lead to intricate interactions among consecutive updates that complicate the analysis of the learning dynamics. Compared to the existing works (Ahn et al., 2023; Zhang et al., 2023; Schlag et al., 2021) that mainly focus on linear self-attention with $\varphi(\cdot)$ being the identify function, our analysis provides more pertinent results for understanding the denoising performance and learning dynamics of attention mechanisms, capturing the nonlinear interactions and transformations across the layers of modern transformer architectures.

378 4 EXPERIMENTAL RESULTS

380

381

382

384

385

386 387

388 389

390

391

392

393 394 395

405

406

In this section, we evaluate our proposed *attention-only transformer* (AoT) architecture on both language and vision tasks. Due to limited computing and engineering resources, the goal of our experimentation is not to outperform state-of-the-art transformers but to verify that AoT can achieve similar or comparable performance on complex language and vision tasks. Hence we believe, while offering a fully interpretable architecture with a layerwise performance guarantee, AoT holds great potential in practical applicability with further engineering development in the future. In all our implementations, we set the operator $\varphi(\cdot)$ in Eq. (3) to be the softmax function.

4.1 DECODER-ONLY TRANSFORMER FOR LANGUAGE TASKS

To study the performance of our architecture on language tasks, we consider the widely used Generative Pre-Training (GPT) task (Radford et al., 2019). In the context of causal language modeling, the goal is to do the next token prediction in a sequence. To adapt to this task, we modify the AoT architecture by changing the MSSA operator to be a causally masked MSSA, i.e., replacing (5) by

$$\boldsymbol{Z}^{(l+1)} = \boldsymbol{Z}^{(l)} + \eta \sum_{k=1}^{K} \boldsymbol{U}_{k}^{(l)} \boldsymbol{U}_{k}^{(l)^{T}} \boldsymbol{Z}^{(l)} \varphi \left(\operatorname{Mask} \left(\boldsymbol{Z}^{(l)^{T}} \boldsymbol{U}_{k}^{(l)} \mathcal{P} \left(\boldsymbol{U}_{k}^{(l)^{T}} \boldsymbol{Z}^{(l)} \right) \right) \right),$$

where $[Mask(A)]_{ij} = a_{ij}$ if $i \leq j$ and $[Mask(A)]_{ij} = -\infty$ otherwise. Following the implement 397 tation used in Kitaev et al. (2020), we apply normalization to the "query matrix" $U_k^{(l)^T} Z^{(l)}$, where 398 $A' = \mathcal{P}(A)$ project each column of $A = [a_1, \dots, a_n] \in \mathbb{R}^{d \times n}$ onto unit sphere, i.e., a' = a/||a||. 399 We follow the same pre-processing and post-processing steps in (Yu et al., 2024, Section 4.1.4). 400 Our implementation of the GPT-2 type transformer and training pipeline is based on the framework 401 outlined in Karpathy (2022).³ In addition, to study the effect of removing the MLP layer, we also 402 train models with MLPs in the first half of transformer blocks, referred as Hybrid, as well as models 403 with MLPs in all blocks, referred as Full MLP. 404

4.1.1 LANGUAGE MODELING

407 **Pre-training language models.** We pre-train AoT-based language models of different sizes and 408 GPT-2 (see Table 1 for model sizes) on OpenWebText (Gokaslan & Cohen, 2019). Here, we train 409 these models over a 1024-token context using the AdamW optimizer (Loshchilov & Hutter, 2019). We plot the training loss and validation loss against the number of training iterations in Figure 5(a) 410 and (b), respectively. It is observed that AoT-based language models of medium and large size can 411 achieve comparable performance to the GPT-2 base model in terms of training and validation loss. In 412 addition, a comparison of AoT models with the Hybrid and Full MLP configurations demonstrates 413 that incorporating MLP layers can accelerate the training process. 414

Zero-shot evaluation. Using the above pre-trained models, we compute the cross-entropy valida-415 tion loss without training on datasets WikiText (Merity et al., 2016)⁴, LAMBADA (Paperno et al., 416 2016)⁵, and PTB (Marcus et al., 1993) in Table 1. In addition, we report zero-shot accuracy in Ta-417 ble 1 on LAMBADA for predicting the final word of sentences, as well as on the Children's Book 418 Test (CBT) (Hill et al., 2016), where the task is to choose either common nouns (CN) or named enti-419 ties (NE) from 10 possible options for an omitted word in a paragraph. It is observed that AoT with 420 medium and large parameter sizes can achieve comparable performance to the GPT-2 base model 421 on these tasks. Moreover, we found that adding MLP layers to AoT does not improve the zero-shot 422 performance. These results highlight the potential of attention-only models to achieve competitive 423 results while maintaining interpretability.

424 425

426

4.1.2 IN-CONTEXT LEARNING ON SIMPLE FUNCTION CLASSES

In-context learning (ICL) refers to the ability of modern language models to perform tasks by using examples provided in the input prompt, along with a new query input, generating outputs without

⁴²⁹ ³https://github.com/karpathy/nanoGPT.git

^{430 &}lt;sup>4</sup>For WikiText2 and WikiText103 (Merity et al., 2016), the test splits are the same, so we merge them as a single dataset referred to as WikiText.

⁵To obtain the accuracy on LAMBADA dataset, we use greedy decoding.





updating the parameters (Brown et al., 2020b; Garg et al., 2023; Park et al., 2024). We evaluate the ICL capabilities of our AoT and compare its performance with that of GPT-2 (Radford et al., 2019). Each model is trained from scratch on specific tasks, including linear and sparse linear regressions. We mainly follow the setup in Garg et al. (2023) to train models to learn linear functions in context. Specifically, for a specific function class \mathcal{G} , we generate random prompts by sampling a function $g \in \mathcal{G}$ from distribution $\mathcal{D}_{\mathcal{G}}$ over functions random inputs $\mathbf{x}_1, \ldots, \mathbf{x}_N \in \mathbb{R}^d$ i.i.d. from $\mathcal{D}_{\mathcal{X}}$ over inputs. To evaluate the inputs on g, we create a prompt $P = (\mathbf{x}_1, g(\mathbf{x}_1), \ldots, \mathbf{x}_N, g(\mathbf{x}_N))$. We train the model f_{θ} to minimize the expected loss over all prompts prefixes:

$$\min_{\boldsymbol{\theta}} \mathbb{E}_{P} \left[\frac{1}{N} \sum_{i=1}^{N-1} \left(f_{\boldsymbol{\theta}}(P^{i}) - g(\boldsymbol{x}_{i}) \right)^{2} \right], \tag{9}$$

where P^i is the prompt prefix up to the input *i*-th in-context example $P = (x_1, q(x_1), \dots, x_i)$.

Tasks. We consider both linear functions and sparse linear functions with dimension d = 20. The in-context examples x_i are sampled from the isotropic Gaussian distribution. For linear functions, we define $\mathcal{G} = \{g : g(x) = w^T x\}$, where x is sampled from the isotropic Gaussian distribution as well. For sparse linear functions, the setup is similar, but with a modification: only 3 coordinates in the vector w are set as non-zero, while the remaining coordinates are set as zero.

Training and evaluation. For all experiments, we set the number of heads to 8 and embedding size 128. To match the sizes of different models by controlling the number of layers. The transformer and Full MLP has 16 layers, Hybrid 24, and AoT 16. To train the model, we sample a batch of random prompts with size 64 and train the models for 50,000 iterations using Adam optimizer (Kingma & Ba, 2017). We evaluate models using same $D_{\mathcal{G}}$ and $D_{\mathcal{X}}$ to sample 1280 prompts. We refer the reader to Park et al. (2024) for more details.

Models # of parameters	$\begin{array}{c c} \textbf{LAMBADA} \\ (val loss) \downarrow \end{array}$	PTB (val loss) ↓	WikiText (val loss) ↓	LAMBADA (acc) ↑	$\begin{array}{c} \textbf{CBT CN} \\ (acc) \uparrow \end{array}$	$\begin{array}{c} \textbf{CBT NE} \\ (acc) \uparrow \end{array}$
Base 102M	4.70	6.03	4.65	0.25	0.80	0.74
Medium 182M	4.47	5.08	4.22	0.29	0.84	0.77
Large 326M	4.26	4.77	3.99	0.34	0.86	0.81
Hybrid 81M	4.84	5.83	4.56	0.25	0.79	0.73
Full MLP 109M	4.73	6.95	4.70	0.30	0.83	0.77
GPT-2 Base 124M	4.32	5.75	4.13	0.40	0.87	0.84

Table 1: Zero-shot results on several benchmark datasets.

We plot the estimation error against in-context samples in Figure 6. It is observed that our AoT architecture can in-context learn linear functions and sparse linear functions, achieving performance close to that of GPT-2 style transformer. Adding MLPs does not improve the in-context learning ability of AoT, which further supports the effectiveness of our attention-only architecture.

4.2 VISION TRANSFORMERS FOR SUPERVISED IMAGE CLASSIFICATION

Now we evaluate the performance of AoT as a backbone architecture for supervised image classification tasks. For further simplification, we do *not* even use LayerNorm layers in the AoT architecture.

1.2 1.2 Transformer Transformer 1.0 AoT 1.0 Least Squares Squared Error Hybrid Error AoT 0.8 0.8 Full MLP Hybrid Full MLP 0.6 Least Squares 0.6 Squared Lasso (alpha=0.01) 0.4 0.40.2 0.2 0.0 0.0 40 0 0 10 20 30 10 20 30 40 In-context Examples In-context Examples (a) Linear regression (b) Sparse linear regression

Figure 6: Evaluating models on in-context learning linear functions. We plot the normalized squared error as a function of in-context examples.

Model architecture. As we mentioned earlier, for vision tasks, we can use an even simpler architecture without the Layernorm (see Figure 3). We use the same pre-processing map and classification head defined in (Yu et al., 2023a, Section 4.1.1) to construct the AoT-based model. Moreover, we consider AoT-based models with different number of parameters and attention layers, as in Table 2.

Table 2: Top-1 accuracy on ImageNet with different models when pre-trained on ImageNet-21K and then fine-tuned on ImageNet-1K.

Models	ImageNet-1K	# of Parameters	# of Layers	
AoT-Base	70.2%	16M	12(Atten)	
AoT-Large	75.7%	52M	24(Atten)	
AoT-Huge	79.2%	86M	32(Atten)	
CRATE- α -B/16 (Yang et al., 2024)	81.2%	72.3M	12(Atten+MLP)	
CRATE- α -L/14 (Yang et al., 2024)	83.9%	253.8M	24 (Atten+MLP)	

Training setup. We employ Lion optimizer (Chen et al., 2024) to pre-train the above AoT-based transformer on ImageNet-21K and AdamW (Loshchilov, 2017) to fine-tune it on ImageNet-1K (Deng et al., 2009) by minimizing the cross-entropy (CE) loss. During the pre-training, we set the learning rate as 1×10^{-4} , weight decay as 0.05, and batch size as 4096. During the fine-tuning, the learning rate as 5×10^{-5} , weight decay as 0.05, and batch size as 2048. Standard data augmentation techniques, including random cropping, random horizontal flipping, and random augmentation, are used in our implementation, the same as those used in Yu et al. (2023b).

⁵²⁴ Based on the above experimental setup, we report the top-1 accuracy of AoT on ImageNet-1K in ⁵²⁵ Table 2. For comparison, we also report the performance of CRATE- α models in Yang et al. (2024), ⁵²⁶ which are enhanced white-box vision models built on CRATE (Yu et al., 2023b). Despite the absence ⁵²⁷ of MLP layers in AoT, it achieves a competitive performance comparable to that of CRATE. This ⁵²⁸ result demonstrates the effectiveness and efficiency of the attention-only architecture.

529 530

531

486

487

488

489

490

491

492

493 494

495

496

497

498

499

500 501 502

503

504

505 506

507

5 CONCLUSION

532 In this work, we propose a new and minimalistic transformer architecture by interpreting each layer 533 as the application of a subspace denoising operator to token representations, where these repre-534 sentations are assumed to be sampled from a mixture of low-rank Gaussians. Remarkably, this 535 architecture consists of subspace self-attention layers and skip connections at each layer, without 536 the MLP operators at all. We have shown that each such layer improves the signal-to-noise ratio 537 of token representations at a linear rate with respect to the number of layers. We have verified the practical potential of this simple architecture through extensive experiments on both language and 538 vision tasks, which strongly suggest that it could lead to more efficient and effective architectures in the future.

540 REFERENCES

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical
 report. *arXiv preprint arXiv:2303.08774*, 2023.
- Kwangjun Ahn, Xiang Cheng, Hadi Daneshmand, and Suvrit Sra. Transformers learn to implement
 preconditioned gradient descent for in-context learning. *arXiv preprint arXiv:2306.00297*, 2023.
- Albert Alcalde, Giovanni Fantuzzi, and Enrique Zuazua. Clustering in pure-attention hardmax trans formers and its role in sentiment analysis. *arXiv preprint arXiv:2407.01602*, 2024.
- Sanjeev Arora, Yuanzhi Li, Yingyu Liang, Tengyu Ma, and Andrej Risteski. Linear algebraic structure of word senses, with applications to polysemy. *Transactions of the Association for Computational Linguistics*, 6:483–495, 2018.
- Fan Bao, Shen Nie, Kaiwen Xue, Yue Cao, Chongxuan Li, Hang Su, and Jun Zhu. All are worth words: A vit backbone for diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 22669–22679, 2023.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In *Advances in Neural Information Processing Systems*, volume 33, pp. 1877–1901, 2020a.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,
 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are
 few-shot learners. Advances in neural information processing systems, 33:1877–1901, 2020b.
- Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and
 Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 9650–9660, 2021.
- Kwan Ho Ryan Chan, Yaodong Yu, Chong You, Haozhi Qi, John Wright, and Yi Ma. Redunet: A white-box deep network from the principle of maximizing rate reduction, 2022.
- Lili Chen, Kevin Lu, Aravind Rajeswaran, Kimin Lee, Aditya Grover, Misha Laskin, Pieter Abbeel,
 Aravind Srinivas, and Igor Mordatch. Decision transformer: Reinforcement learning via sequence
 modeling. Advances in neural information processing systems, 34:15084–15097, 2021.
- Mark Chen, Alec Radford, Rewon Child, Jeffrey Wu, Heewoo Jun, David Luan, and Ilya Sutskever.
 Generative pretraining from pixels. In *International conference on machine learning*, pp. 1691–
 1703. PMLR, 2020.
- Xiangning Chen, Chen Liang, Da Huang, Esteban Real, Kaiyuan Wang, Hieu Pham, Xuanyi Dong, Thang Luong, Cho-Jui Hsieh, Yifeng Lu, et al. Symbolic discovery of optimization algorithms. *Advances in neural information processing systems*, 36, 2024.
- Yubei Chen, Zeyu Yun, Yi Ma, Bruno Olshausen, and Yann LeCun. Minimalistic unsupervised
 representation learning with the sparse manifold transform. In *The Eleventh International Con- ference on Learning Representations*, 2023.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hi erarchical image database. In 2009 IEEE conference on computer vision and pattern recognition,
 pp. 248–255. Ieee, 2009.
- Jacob Devlin. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- Yihe Dong, Jean-Baptiste Cordonnier, and Andreas Loukas. Attention is not all you need: Pure attention loses rank doubly exponentially with depth. In *International Conference on Machine Learning*, pp. 2793–2803. PMLR, 2021.

594 595	Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An
596 597	image is worth 16x16 words: Transformers for image recognition at scale. <i>arXiv preprint arXiv:2010.11929</i> , 2020.
598	
599	Nelson Elhage, Tristan Hume, Catherine Olsson, Nicholas Schiefer, Tom Henighan, Shauna Kravec,
600	Zac Hatfield-Dodds, Robert Lasenby, Dawn Drain, Carol Chen, et al. Toy models of superposi-
601	tion. arxiv preprint arXiv:2209.10652, 2022.
602	Joshua Engels, Isaac Liao, Eric J Michaud, Wes Gurnee, and Max Tegmark. Not all language model
603	features are linear. arXiv preprint arXiv:2405.14860, 2024.
604	
605	Shivam Garg, Dimitris Tsipras, Percy Liang, and Gregory Valiant. What can transformers learn
606 607	in-context? a case study of simple function classes, 2023. URL https://arxiv.org/abs/ 2208.01066.
608	Borian Geshkovski, Cyril Letrouit, Yury Polyanskiy, and Philippe Rigollet. The emergence of clus-
609	ters in self-attention dynamics. <i>arXiv preprint arXiv:2305.05465</i> , 2023a.
611 612	Borjan Geshkovski, Cyril Letrouit, Yury Polyanskiy, and Philippe Rigollet. A mathematical per- spective on transformers. <i>arXiv preprint arXiv:2312.10794</i> , 2023b.
613 614	Mor Geva, Roei Schuster, Jonathan Berant, and Omer Levy. Transformer feed-forward layers are key-value memories. <i>arXiv preprint arXiv:2012.14913</i> , 2020.
615 616	Aaron Gokaslan and Vanya Cohen. Openwebtext corpus. http://Skylion007.github.io/
617	openwebrexccorpus, 2019.
618	Karol Gregor and Yann LeCun. Learning fast approximations of sparse coding. In Proceedings of
619	the 27th international conference on international conference on machine learning, pp. 399–406,
620	2010.
621	Jialong Guo, Xinghao Chen, Yehui Tang, and Yunhe Wang. Slah: Efficient transformers with
622 623	simplified linear attention and progressive re-parameterized batch normalization. <i>arXiv preprint arXiv:2405.11582</i> , 2024.
624	Dathy Us and Thomas Hafmann Simplifying transformer blocks. In The Twelfth International
625 626	Conference on Learning Representations, 2024.
627 628 629	Felix Hill, Antoine Bordes, Sumit Chopra, and Jason Weston. The goldilocks principle: Reading children's books with explicit memory representations, 2016. URL https://arxiv.org/abs/1511.02301.
631 632	Yibo Jiang, Goutham Rajendran, Pradeep Ravikumar, Bryon Aragam, and Victor Veitch. On the origins of linear representations in large language models. <i>arXiv preprint arXiv:2403.03867</i> , 2024.
634	Andrej Karpathy. NanoGPT. https://github.com/karpathy/nanoGPT, 2022.
635	Angelos Katharonoulos Anoory Vyes Nikoleos Dennes and Erenaois Elevent Transformers are
030	rnns: Fast autoregressive transformers with linear attention. In International conference on ma-
638	chine learning, pp. 5156–5165. PMLR, 2020.
639	Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. 2017. URL
640 641	https://arxiv.org/abs/1412.6980.
642	Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. <i>arXiv</i>
643	preprint arXiv:2304.02643, 2023.
645	Nilita Kitaan Lukaan Kaisan and Analys Laurians. Deferming The Solid tangles V
646	INIKITA KITACV, LUKASZ KAISEF, AND ANSEIM LEVSKAYA. KETORMET: THE EINCIENT TRANSFORMET. <i>arXiv</i>
647	preprint arAiv.2001.04451, 2020.

I Loshchilov. Decoupled weight decay regularization. arXiv preprint arXiv:1711.05101, 2017.

657

659

685

687 688

689

- 648 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization, 2019. URL https: 649 //arxiv.org/abs/1711.05101. 650
- Jingi Luo, Tianjiao Ding, Kwan Ho Ryan Chan, Darshan Thaker, Aditya Chattopadhyay, Chris 651 Callison-Burch, and René Vidal. Pace: Parsimonious concept engineering for large language 652 models. arXiv preprint arXiv:2406.04331, 2024. 653
- 654 Yi Ma, Doris Tsao, and Heung-Yeung Shum. On the principles of parsimony and self-consistency 655 for the emergence of intelligence. Frontiers of Information Technology & Electronic Engineering, 656 23(9):1298-1323, 2022.
- Mitchell P. Marcus, Beatrice Santorini, and Mary Ann Marcinkiewicz. Building a large annotated 658 corpus of English: The Penn Treebank. Computational Linguistics, 19(2):313–330, 1993. URL https://aclanthology.org/J93-2004. 660
- 661 Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mixture models, 2016. 662
- 663 Vishal Monga, Yuelong Li, and Yonina C Eldar. Algorithm unrolling: Interpretable, efficient deep 664 learning for signal and image processing. *IEEE Signal Processing Magazine*, 38(2):18–44, 2021. 665
- 666 Lorenzo Noci, Chuning Li, Mufan Li, Bobby He, Thomas Hofmann, Chris J Maddison, and Dan 667 Roy. The shaped transformer: Attention models in the infinite depth-and-width limit. Advances in Neural Information Processing Systems, 36, 2024. 668
- 669 Druv Pai, Sam Buchanan, Ziyang Wu, Yaodong Yu, and Yi Ma. Masked completion via structured 670 diffusion with white-box transformers. In The Twelfth International Conference on Learning 671 Representations, 2023. 672
- Denis Paperno, Germán Kruszewski, Angeliki Lazaridou, Quan Ngoc Pham, Raffaella Bernardi, 673 Sandro Pezzelle, Marco Baroni, Gemma Boleda, and Raquel Fernández. The lambada dataset: 674 Word prediction requiring a broad discourse context, 2016. URL https://arxiv.org/ 675 abs/1606.06031. 676
- 677 Jongho Park, Jaeseung Park, Zheyang Xiong, Nayoung Lee, Jaewoong Cho, Samet Oymak, Kang-678 wook Lee, and Dimitris Papailiopoulos. Can mamba learn how to learn? a comparative study on 679 in-context learning tasks, 2024. URL https://arxiv.org/abs/2402.04248.
- 680 Kiho Park, Yo Joong Choe, and Victor Veitch. The linear representation hypothesis and the geometry 681 of large language models. arXiv preprint arXiv:2311.03658, 2023. 682
- 683 William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of* 684 the IEEE/CVF International Conference on Computer Vision, pp. 4195–4205, 2023.
- Telmo Pessoa Pires, António V Lopes, Yannick Assogba, and Hendra Setiawan. One wide feedfor-686 ward is all you need. arXiv preprint arXiv:2309.01826, 2023.
 - Phil Pope, Chen Zhu, Ahmed Abdelkader, Micah Goldblum, and Tom Goldstein. The intrinsic dimension of images and its impact on learning. In International Conference on Learning Representations, 2020.
- 691 Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language 692 models are unsupervised multitask learners. OpenAI blog, 1(8):9, 2019. 693
- Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar 694 Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. Advances in Neural Informa-696 tion Processing Systems, 35:36479–36494, 2022. 697
- 698 Imanol Schlag, Kazuki Irie, and Jürgen Schmidhuber. Linear transformers are secretly fast weight 699 programmers. In International Conference on Machine Learning, pp. 9355–9366. PMLR, 2021. 700
- Sainbayar Sukhbaatar, Edouard Grave, Guillaume Lample, Herve Jegou, and Armand Joulin. Aug-701 menting self-attention with persistent memory. arXiv preprint arXiv:1907.01470, 2019.

- Xiaoxia Sun, Nasser M Nasrabadi, and Trac D Tran. Supervised deep sparse coding networks for image classification. *IEEE Transactions on Image Processing*, 29:405–418, 2019.
- Adly Templeton. Scaling monosemanticity: Extracting interpretable features from claude 3 sonnet.
 Anthropic, 2024.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation language models, 2023. URL https://arxiv.org/abs/2302.13971.
- Yao-Hung Hubert Tsai, Shaojie Bai, Makoto Yamada, Louis-Philippe Morency, and Ruslan Salakhutdinov. Transformer dissection: An unified understanding for transformer's attention via the lens of kernel. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pp. 4344–4353, 2019.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. Advances in neural information processing systems, 30, 2017.
- Roman Vershynin. *High-dimensional probability: An introduction with applications in data science*, volume 47. Cambridge university press, 2018.
- Shenlong Wang, Sanja Fidler, and Raquel Urtasun. Proximal deep structured models. Advances in Neural Information Processing Systems, 29, 2016.
- John Wright and Yi Ma. *High-dimensional data analysis with low-dimensional models: Principles, computation, and applications.* Cambridge University Press, 2022.
- Xinyi Wu, Amir Ajorlou, Yifei Wang, Stefanie Jegelka, and Ali Jadbabaie. On the role of attention
 masks and layernorm in transformers. *arXiv preprint arXiv:2405.18781*, 2024.
- Jinrui Yang, Xianhang Li, Druv Pai, Yuyin Zhou, Yi Ma, Yaodong Yu, and Cihang Xie. Scaling white-box transformers for vision. *arXiv preprint arXiv:2405.20299*, 2024.
- Yaodong Yu, Sam Buchanan, Druv Pai, Tianzhe Chu, Ziyang Wu, Shengbang Tong, Hao Bai, Yuexiang Zhai, Benjamin D Haeffele, and Yi Ma. White-box transformers via sparse rate reduction: Compression is all there is? *arXiv preprint arXiv:2311.13110*, 2023a.
- Yaodong Yu, Sam Buchanan, Druv Pai, Tianzhe Chu, Ziyang Wu, Shengbang Tong, Benjamin D Haeffele, and Yi Ma. White-box transformers via sparse rate reduction. *arXiv preprint arXiv:2306.01129*, 2023b.
- Yaodong Yu, Sam Buchanan, Druv Pai, Tianzhe Chu, Ziyang Wu, Shengbang Tong, Hao Bai, Yuexiang Zhai, Benjamin D. Haeffele, and Yi Ma. White-box transformers via sparse rate reduction: Compression is all there is?, 2024. URL https://arxiv.org/abs/2311.13110.
- Zeyu Yun, Yubei Chen, Bruno A Olshausen, and Yann LeCun. Transformer visualization via dictio nary learning: contextualized embedding as a linear superposition of transformer factors. *arXiv preprint arXiv:2103.15949*, 2021.
- Jian Zhang and Bernard Ghanem. Ista-net: Interpretable optimization-inspired deep network for image compressive sensing. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 1828–1837, 2018.
- Ruiqi Zhang, Spencer Frei, and Peter L Bartlett. Trained transformers learn linear models in-context.
 arXiv preprint arXiv:2306.09927, 2023.
- 752
- 753
- 754
- 755

To simplify our development, we introduce some further notation. We use BlkDiag (X_1, \ldots, X_K) to denote a block diagonal matrix whose diagonal blocks are X_1, \ldots, X_K .

A PROOF OF THEOREM 1

A.1 PRELIMINARY RESULTS

To prove Theorem 1, we first establish several probabilistic results about Gaussian random vectors. First, we present a probabilistic bound on the deviation of the norm of Gaussian random vectors from its mean. This is an extension of (Vershynin, 2018, Theorem 3.1.1).

Lemma 1. Let $x \sim \mathcal{N}(\mathbf{0}, \delta^2 \mathbf{I}_d)$ be a Gaussian random vector. It holds with probability at least $1 - 2 \exp(-t^2/2\delta^2)$ that

$$\left| \|\boldsymbol{x}\| - \delta \sqrt{d} \right| \le t + 2\delta. \tag{10}$$

Based on the above lemma, we can respectively estimate the norm of coefficients in the signal and noise parts, the products between different pairs of Gaussian random vectors, and the bounds on the soft-max values of these products.

Lemma 2. Consider the setting in Definition 1 with $p = p_1 = \cdots = p_K$ and $N_1 = \cdots = N_K = N/K$. Suppose that $p \ge 16(\sqrt{\log N} + 1)^2$ and

$$N \ge 8\pi K^2 \log^3 N, \ \delta \le \frac{1}{8} \sqrt{\frac{\log N}{p}}.$$
(11)

The following statements hold:

(*i*) With probability at least $1 - 2KN^{-1}$, we have

$$|\|\boldsymbol{a}_i\| - \sqrt{p}| \le 2\left(\sqrt{\log N} + 1\right), \forall i \in [N],$$
(12)

$$|\|\boldsymbol{e}_{i,l}\| - \delta\sqrt{p}| \le 2\delta\left(\sqrt{\log N} + 1\right), \forall i \in C_k, l \neq k \in [K].$$
(13)

(*ii*) With probability at least $1 - 4KN^{-2}$, we have

$$\langle \boldsymbol{a}_i, \boldsymbol{a}_j \rangle | \le 3\sqrt{\log N} \| \boldsymbol{a}_i \|, \forall i \neq j \in C_k, k \in [K],$$
(14)

$$|\langle \boldsymbol{a}_i, \boldsymbol{e}_{j,l} \rangle| \le 3\sqrt{\log N} \|\boldsymbol{e}_{j,l}\|, \forall i \in C_k, j \in C_l, k \neq l \in [K],$$
(15)

$$|\langle \boldsymbol{e}_{i,k}, \boldsymbol{e}_{j,k} \rangle| \le 3\delta \sqrt{\log N} \|\boldsymbol{e}_{j,k}\|, \forall i \in C_l, j \in C_m, l, m \neq k.$$
(16)

(iii) With probability at least $1 - 2N^{-1}$, we have

$$\max_{i \in C_k} \langle \boldsymbol{a}_i, \boldsymbol{e}_{j,k} \rangle \ge \sqrt{\log N} \| \boldsymbol{e}_{j,k} \|, \forall j \in C_l, l \neq k \in [K].$$
(17)

(iv) With probability at least $1 - 4KN^{-1}$, we have

$$\frac{\exp\left(\langle \boldsymbol{a}_{i}, \boldsymbol{e}_{j,k} \rangle\right)}{\sum_{i' \in C_{k}} \exp\left(\langle \boldsymbol{a}_{i'}, \boldsymbol{e}_{j,k} \rangle\right)} \leq \frac{1}{2}, \forall i \in C_{k}, j \in C_{l}, k \neq l \in [K],$$
(18)

$$\frac{\exp\left(\langle \boldsymbol{e}_{i,k}, \boldsymbol{e}_{j,k} \rangle\right)}{\sum_{i' \neq j, i' \in C_l} \exp\left(\langle \boldsymbol{e}_{i',k}, \boldsymbol{e}_{j,k} \rangle\right)} \le \frac{1}{2}, \forall i \neq j, i \in C_l, j \in C_m, l, m \neq k.$$
(19)

Proof. (i) Applying Lemma 1 to $a_i \sim \mathcal{N}(\mathbf{0}, I_p)$ with $t = 2\sqrt{\log N}$ yields

$$\mathbb{P}\left(\left|\left|\left|\boldsymbol{a}_{i}\right|\right|-\sqrt{p}\right| \leq 2(\sqrt{\log N}+1)\right) \geq 1-2N^{-2}.$$

This, together with the union bound, yields that (12) holds for all $i \in [N]$ with probability at least $1-2N^{-1}$. Using the same argument, we obtain that (13) holds for all $i \in C_k$ and $l \neq k \in [K]$ with probability at least $1-2(K-1)N^{-1}$. Finally, applying the union bound yields that the probability is $1-2KN^{-1}$.

(ii) For each pair (i, j) with $i \neq j \in C_k$ and $k \in [K]$, conditioned on a_i , we have $\langle a_i, a_j \rangle \sim \mathcal{N}(0, \|a_i\|^2)$. According to the tail bound the Gaussian random variable, we have

$$\mathbb{P}\left(|\langle \boldsymbol{a}_i, \boldsymbol{a}_j\rangle| \geq 3 \|\boldsymbol{a}_i\| \sqrt{\log N} |\boldsymbol{a}_i\right) \leq 2N^{-4}.$$

This, together with the union bound, implies that conditioned on a_i , it holds with probability at least $1 - 2N^{-2}$ that $|\langle a_i, a_j \rangle| \le 2 ||a_i|| \sqrt{\log N}$ for all $i \ne j \in C_k$ and $k \in [K]$. Using the same argument, we obtain (15) and (16). Finally, applying the union bound yields the probability.

(iii) Conditioned on $e_{j,k}$, we obtain that $X_i := \langle a_i, e_{j,k} \rangle / ||e_{j,k}|| \sim \mathcal{N}(0,1)$ for each $i \in C_k$ are i.i.d. standard normal random variables. Then, we have

$$\mathbb{P}\left(\max_{i\in C_k} X_i \ge \sqrt{\log N}\right) = 1 - \left(\mathbb{P}\left(X_1 < \sqrt{\log N}\right)\right)^{N_k}.$$
(20)

Using the property of the standard Gaussian random variable, we have

$$\mathbb{P}(X_1 \ge t) \ge \left(\frac{1}{t} - \frac{1}{t^3}\right) \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{t^2}{2}\right)$$

Taking $t = \sqrt{\log N}$, we obtain

$$\mathbb{P}\left(X_1 \ge \sqrt{\log N}\right) = \frac{1}{\sqrt{\log N}} \left(1 - \frac{1}{\log N}\right) \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{\log N}{2}\right) \ge \frac{1}{2\sqrt{2\pi N \log N}},$$
 (21)

where the inequality follows from $N \ge \exp(2)$. Substituting this into (20) yields

$$\mathbb{P}\left(\max_{i\in C_{k}}X_{i}\geq\sqrt{\log N}\right)\geq1-\left(1-\frac{1}{2\sqrt{2\pi N\log N}}\right)^{N/K}$$
$$\geq1-\exp\left(-\frac{\sqrt{N}}{2K\sqrt{2\pi\log N}}\right)\geq1-N^{-1},$$

where the second inequality uses $1 - x \le \exp(-x)$ for all x > 0 and the last inequality follows from $N \ge 8\pi K^2 \log^3 N$. This, together with the definition of X_i , implies (17).

(iv) Conditioned on $e_{j,k}$, we have $X_i := \langle a_i, e_{j,k} \rangle \sim \mathcal{N}(0, ||e_{j,k}||^2)$ for each $i \in C_k$ are i.i.d. normal random variables. Suppose that (13) holds for all $i \in C_k, l \neq k \in [K]$, which happens with probability at least $1 - 2(K - 1)N^{-1}$ according to (i). This implies for all $j \in C_k$ and $k \in [K]$,

$$\|\boldsymbol{e}_{j,k}\| \le \delta\left(\sqrt{p} + 2\sqrt{\log N} + 2\right) \le \frac{3}{2}\delta\sqrt{p},\tag{22}$$

where the last inequality follows from $p \ge 16(\sqrt{\log N} + 1)^2$ due to (11). For ease of exposition, let

$$\sigma := \| \boldsymbol{e}_{j,k} \|, \ S := \sum_{i \in C_k} \exp(X_i).$$
(23)

Obviously, showing (18) is equivalent to proving

$$2\exp(X_i) \le \sum_{i' \in C_k} \exp\left(X_{i'}\right) = S, \ \forall i \in C_k.$$

$$(24)$$

Note that $X_i/\sigma \sim \mathcal{N}(0,1)$ for all $i \in C_k$. Using the tail bound of the standard normal random variable, we have

$$\mathbb{P}\left(\frac{|X_i|}{\sigma} \ge 2\sqrt{\log N}\right) \le 2N^{-2}, \ \forall i \in C_k.$$

This, together with the union bound, yields that it holds with probability $1 - 2N^{-1}$ that $|X_i| \le 2\sigma\sqrt{\log N}$ for all $i \in [N]$. Using this, (22), (23), and the union bound, we obtain with probability at least $1 - 2KN^{-1}$,

$$|X_i| \le 3\delta\sqrt{p\log N}, \ \forall i \in [N].$$

Therefore, we have 865

$$\exp\left(-3\delta\sqrt{p\log N}\right) \le \exp(X_i) \le \exp\left(3\delta\sqrt{p\log N}\right), \ \forall i \in [N].$$
(25)

Using this and (23), we have

$$S \geq \frac{N}{K} \exp\left(-3\delta \sqrt{p \log N}\right)$$

872 This, together with (25), implies that proving (24) is sufficient to proving

 $\log N \ge 6\delta \sqrt{p \log N} + \log \left(2K\right),$

which holds when $N \ge \max\{16K^4, \exp(64\delta^2 p)\}$ due to (11). According to the union bound, (18) holds with probability at least $1 - 2KN^{-1}$. Using the same argument, (19) holds with probability at least $1 - 2KN^{-1}$.

 $M_{1} := \begin{bmatrix} \theta^{2} \boldsymbol{A}_{1}^{T} \boldsymbol{A}_{1} & \theta \boldsymbol{A}_{1}^{T} \boldsymbol{E}_{2,1} & \dots & \theta \boldsymbol{A}_{1}^{T} \boldsymbol{E}_{K,1} \\ \theta \boldsymbol{E}_{2,1}^{T} \boldsymbol{A}_{1} & \boldsymbol{E}_{2,1}^{T} \boldsymbol{E}_{2,1} & \dots & \boldsymbol{E}_{2,1}^{T} \boldsymbol{E}_{K,1} \\ \vdots & \vdots & \ddots & \vdots \\ \theta \boldsymbol{E}_{K,1}^{T} \boldsymbol{A}_{1} & \boldsymbol{E}_{K,1}^{T} \boldsymbol{E}_{2,1} & \dots & \boldsymbol{E}_{K,1}^{T} \boldsymbol{E}_{K,1} \end{bmatrix} \in \mathbb{R}^{N \times N},$ $M_{2} := \begin{bmatrix} \boldsymbol{E}_{1,2}^{T} \boldsymbol{E}_{1,2} & \theta \boldsymbol{E}_{1,2}^{T} \boldsymbol{A}_{2} & \dots & \boldsymbol{E}_{1,2}^{T} \boldsymbol{E}_{K,2} \\ \theta \boldsymbol{A}_{2}^{T} \boldsymbol{E}_{1,2}^{T} & \theta^{2} \boldsymbol{A}_{2}^{T} \boldsymbol{A}_{2} & \dots & \theta \boldsymbol{A}_{2}^{T} \boldsymbol{E}_{K,2} \\ \vdots & \vdots & \ddots & \vdots \\ \boldsymbol{E}_{K,2}^{T} \boldsymbol{E}_{1,2} & \theta \boldsymbol{E}_{K,2}^{T} \boldsymbol{A}_{2} & \dots & \boldsymbol{E}_{K,2}^{T} \boldsymbol{E}_{K,2} \end{bmatrix} \in \mathbb{R}^{N \times N},$

A.2 PROOF OF THEOREM 1

To simplify our development, let

where $\theta \geq 1$. Recall that

$$\boldsymbol{Z}^{(0)} = \begin{bmatrix} \boldsymbol{Z}_1^{(0)} & \dots & \boldsymbol{Z}_K^{(0)} \end{bmatrix} = \begin{bmatrix} \boldsymbol{U}_1 \boldsymbol{A}_1 + \sum_{j \neq 1} \boldsymbol{U}_j \boldsymbol{E}_{1,j} & \dots & \boldsymbol{U}_K \boldsymbol{A}_K + \sum_{j \neq K} \boldsymbol{U}_j \boldsymbol{E}_{K,j} \end{bmatrix},$$
(27)

 $\boldsymbol{M}_{K} := \begin{bmatrix} \boldsymbol{E}_{1,K}^{T} \boldsymbol{E}_{1,K} & \boldsymbol{E}_{1,K}^{T} \boldsymbol{E}_{2,K} & \dots & \boldsymbol{\theta} \boldsymbol{E}_{1,K}^{T} \boldsymbol{A}_{K} \\ \boldsymbol{E}_{2,K}^{T} \boldsymbol{E}_{1,K} & \boldsymbol{E}_{2,K}^{T} \boldsymbol{E}_{2,K} & \dots & \boldsymbol{\theta} \boldsymbol{E}_{2,K}^{T} \boldsymbol{A}_{K} \\ \vdots & \vdots & \ddots & \vdots \\ \boldsymbol{\theta} \boldsymbol{A}^{T} \boldsymbol{E} & \boldsymbol{\theta} \boldsymbol{A}^{T} \boldsymbol{E} & \boldsymbol{\theta} \boldsymbol{A}^{T} \boldsymbol{E} & \boldsymbol{\theta}^{T} \boldsymbol{A} \end{bmatrix} \in \mathbb{R}^{N \times N}.$

Lemma 3. Consider the setting in Definition 1 with $p = p_1 = \cdots = p_K$ and $N_1 = \cdots = N_K = N/K$. Let $\varphi(\cdot)$ be

$$\varphi(\boldsymbol{x}) = h(\sigma(\boldsymbol{x})), \tag{28}$$

(26)

910 where $\sigma : \mathbb{R}^N \to \mathbb{R}^N$ is the soft-max function and $h : \mathbb{R}^N \to \mathbb{R}^N$ is an element-wise thresholding 911 function with $h(x) = \tau \mathbb{I} \{x > \tau\}$ for each $i \in [N]$. Suppose that (11) holds. Suppose in addition 912 that $p \ge 64(\sqrt{\log N} + 1)^2$ and

$$au \in \left(\frac{1}{2}, \frac{1}{1+N\exp(-9p/32)}\right]$$
(29)

⁹¹⁶ The following statements hold with probability at least $1 - KN^{-\Omega(1)}$ that,

$$\varphi(\boldsymbol{M}_1) = \text{BlkDiag}(\tau \boldsymbol{I}, \boldsymbol{0}, \dots, \boldsymbol{0}), \dots, \varphi(\boldsymbol{M}_K) = \text{BlkDiag}(\boldsymbol{0}, \boldsymbol{0}, \dots, \tau \boldsymbol{I}).$$
(30)

Proof. Suppose that (12)-(19) hold, which happens with probability at least $1 - KN^{-\Omega(1)}$ according to Lemma 2, (11), and the union bound. Now, we focus on studying M_1 as defined in (26). For ease of exposition, we denote the *i*-th column of M_1 by $m_i \in \mathbb{R}^N$ for each $i \in [N]$. Moreover, recall that

$$C_1 = \left\{1, 2, \dots, \frac{N}{K}\right\}, \dots, C_K = \left\{\frac{(K-1)N}{K} + 1, \frac{(K-1)N}{K} + 2, \dots, N\right\}.$$

We now divide our proof into two cases. We first study the *i*-th column of M_1 for each $i \in C_1$, and then study the *i*-th column of M_1 for each $i \in C_k$ with $k \neq 1$.

Case 1. According to (26), we have for each $i \in C_1$,

$$m_{ij} = \theta^2 \langle \boldsymbol{a}_i, \boldsymbol{a}_j \rangle, \forall j \in C_1, \ m_{ij} = \theta \langle \boldsymbol{a}_i, \boldsymbol{e}_{j,k} \rangle, \forall j \in C_k, k \neq 1.$$

For each pair (i, j) with $i \neq j \in C_1$, we compute

$$\frac{\sigma_i(\boldsymbol{m}_i)}{\sigma_j(\boldsymbol{m}_i)} = \exp\left(m_{ii} - m_{ij}\right) \ge \exp\left(\theta \|\boldsymbol{a}_i\| \left(\theta \|\boldsymbol{a}_i\| - 3\sqrt{\log N}\right)\right) \ge \exp\left(\frac{9\theta^2 p}{32}\right), \quad (31)$$

where the first inequality follows from (14) and the second uses (12) and $\sqrt{p} \ge 8(\sqrt{\log N} + 1)$. Using the same argument, for each pair (i, j) with $i \in C_1, j \in C_k$, and $k \neq 1$, we obtain

$$\frac{\sigma_i(\boldsymbol{m}_i)}{\sigma_j(\boldsymbol{m}_i)} \ge \exp\left(\frac{9\theta^2 p}{32}\right)$$

This, together with $\sum_{j=1}^{N} \sigma_j(\boldsymbol{m}_i) = 1$, yields $\left(1 + (N-1) \exp\left(-9\theta^2 p/32\right)\right) \sigma_i(\boldsymbol{m}_i) \ge 1$. There-fore, we have for each $i \in C_1$,

$$\sigma_i(\boldsymbol{m}_i) \ge \frac{1}{1 + N \exp(-9\theta^2 p/32)} > \frac{1}{2}, \ \sigma_j(\boldsymbol{m}_i) \le \frac{1}{2}, \ \forall j \ne i,$$
(32)

where the last inequality follows from $p \ge 64(\sqrt{\log N} + 1)^2$. This, together with the value of τ in (29), yields for each $i \in C_1$,

$$\sigma_j(\boldsymbol{m}_i) < \tau < \sigma_i(\boldsymbol{m}_i), \ \forall j \neq i.$$

Using this and (28), we have for each $i \in C_1$,

$$h(\sigma_i(\boldsymbol{m}_i)) = \tau, \ h(\sigma_j(\boldsymbol{m}_i)) = 0, \ \forall j \neq i.$$

Case 2. For each $i \in C_k$ with $k \neq 1$, it follows from (26) that

$$m_{ij} = \theta \langle \boldsymbol{e}_{i,1}, \boldsymbol{a}_j \rangle, \forall j \in C_1, \ m_{ij} = \langle \boldsymbol{e}_{i,1}, \boldsymbol{e}_{j,1} \rangle, \ \forall j \in C_l, l \neq 1$$

Consider a fixed $i \in C_k$ with $k \neq 1$, it follows from (17) that there exists $j_i \in C_1$ such that $m_{ij_i} \geq \theta \| \boldsymbol{e}_{i,1} \| \sqrt{\log N}$. This implies

$$\frac{\sigma_{j_i}(\boldsymbol{m}_i)}{\sigma_i(\boldsymbol{m}_i)} = \exp\left(\theta m_{ij_i} - m_{ii}\right) \ge \exp\left(\|\boldsymbol{e}_{i,1}\| \left(\theta\sqrt{\log N} - \|\boldsymbol{e}_{i,1}\|\right)\right)$$
$$\ge \exp\left(\frac{3\delta\theta}{4}\sqrt{p\log N} - \frac{25}{16}\delta^2 p\right),$$

where the second inequality follows from (13). This, together with $\sigma_i(\boldsymbol{m}_i) + \sigma_{ii}(\boldsymbol{m}_i) < 1$, implies

$$\sigma_i(\boldsymbol{m}_i) < \frac{1}{1 + \exp\left(3\delta\theta\sqrt{p\log N}/4 - 25\delta^2 p/16\right)} < \frac{1}{1 + \exp\left(\delta\theta\sqrt{p\log N}/2\right)} < \frac{1}{2}, \quad (33)$$

where the second inequality uses $\delta \sqrt{p} \leq \sqrt{\log N}/8$ due to (11). On the other hand, it follows from (18) and (19) that

971
$$\sigma_j(\boldsymbol{m}_i) \leq rac{1}{2}, \forall j \neq i.$$

This, together with (33), $\delta \leq 1/8$, $\sqrt{p} \geq 8(\sqrt{\log N} + 1)$, and the value of τ by (29), yields for each $i \in C_k$ with $k \neq 1$,

$$\sigma_j(\boldsymbol{m}_i) < \tau, \forall j \in [N]. \tag{34}$$

This directly implies

× 7

$$h(\sigma(\boldsymbol{m}_i)) = \mathbf{0}, \ \forall i \in C_k, k \neq 1$$

Then, we have $\varphi(M_1) = \begin{bmatrix} \tau I & 0 \\ 0 & 0 \end{bmatrix}$. Applying the same argument to M_2, \ldots, M_K , we obtain (30).

Armed with the above result, we are ready to prove Theorem 1.

Proof of Theorem 1. For ease of exposition, let $M_k^{(l)} := Z^{(l)^T} U_k U_k^T Z^{(l)}$ for each $k \in [K]$ and $l \in [L]$. Suppose that (30) holds, which happens with probability at least $1 - KN^{-\Omega(1)}$ according to (11), and (29), Lemma 3. We claim that for each $l \in [L]$, we have

$$\boldsymbol{Z}^{(l)} = \begin{bmatrix} (1+\eta\tau)^{l} \boldsymbol{U}_{1}\boldsymbol{A}_{1} + \sum_{j\neq 1} \boldsymbol{U}_{j}\boldsymbol{E}_{1,j} & \dots & (1+\eta\tau)^{l} \boldsymbol{U}_{K}\boldsymbol{A}_{K} + \sum_{j\neq K} \boldsymbol{U}_{j}\boldsymbol{E}_{K,j} \end{bmatrix}.$$
 (35)

This, together with (6), yields for each $k \in [K]$ and $l \in [L]$,

$$\operatorname{SNR}(\boldsymbol{Z}_{k}^{(l)}) = \frac{\|\boldsymbol{U}_{k}\boldsymbol{U}_{k}^{T}\boldsymbol{Z}_{k}^{(l)}\|_{F}}{\|(\boldsymbol{I} - \boldsymbol{U}_{k}\boldsymbol{U}_{k}^{T})\boldsymbol{Z}_{k}^{(l)}\|_{F}} = \frac{(1 + \eta\tau)^{l}\|\boldsymbol{A}_{k}\|_{F}}{\|\sum_{j \neq k}\boldsymbol{U}_{j}\boldsymbol{E}_{k,j}\|_{F}},$$

which directly implies (8) for each $k \in [K]$ and $l \in [L-1]$. According to the union bound, the probability is $1 - KLN^{-\Omega(1)}$.

The rest of the proof is devoted to proving the claim (35) using the induction method. First, we consider the base case l = 1. According to (27) and (7), we compute

$$\boldsymbol{U}_{1}\boldsymbol{U}_{1}^{T}\boldsymbol{Z}^{(0)} = \begin{bmatrix} \boldsymbol{U}_{1}\boldsymbol{A}_{1} & \boldsymbol{U}_{1}\boldsymbol{E}_{2,1} & \dots & \boldsymbol{U}_{1}\boldsymbol{E}_{K,1} \end{bmatrix}, \\ \boldsymbol{M}_{1}^{(0)} = (\boldsymbol{U}_{1}\boldsymbol{U}_{1}^{T}\boldsymbol{Z}^{(0)})^{T}(\boldsymbol{U}_{1}\boldsymbol{U}_{1}^{T}\boldsymbol{Z}^{(0)}) = \begin{bmatrix} \boldsymbol{A}_{1}^{T}\boldsymbol{A}_{1} & \boldsymbol{A}_{1}^{T}\boldsymbol{E}_{2,1} & \dots & \boldsymbol{A}_{1}^{T}\boldsymbol{E}_{K,1} \\ \boldsymbol{E}_{2,1}^{T}\boldsymbol{A}_{1} & \boldsymbol{E}_{2,1}^{T}\boldsymbol{E}_{2,1} & \dots & \boldsymbol{E}_{2,1}^{T}\boldsymbol{E}_{K,1} \\ \vdots & \vdots & \ddots & \vdots \\ \boldsymbol{E}_{K,1}^{T}\boldsymbol{A}_{1} & \boldsymbol{E}_{K,1}^{T}\boldsymbol{E}_{2,1} & \dots & \boldsymbol{E}_{K,1}^{T}\boldsymbol{E}_{K,1} \end{bmatrix}$$

Using the same argument, we can compute $M_k^{(0)}$ for each $k \in [K]$. This, together with (30) for each $k \in [K]$, yields

$$\sum_{k=1}^{K} \boldsymbol{U}_{k} \boldsymbol{U}_{k}^{T} \boldsymbol{Z}^{(0)} \varphi(\boldsymbol{M}_{k}^{(0)}) = \begin{bmatrix} \tau \boldsymbol{U}_{1} \boldsymbol{A}_{1} & \tau \boldsymbol{U}_{2} \boldsymbol{A}_{2} & \dots & \tau \boldsymbol{U}_{K} \boldsymbol{A}_{K} \end{bmatrix}.$$

Using this, (27), and (3), we directly obtain that (35) holds for l = 1. Next, we consider the case $l \ge 2$. Suppose that (35) holds for some $l \ge 1$. We compute

$$\boldsymbol{U}_{1}\boldsymbol{U}_{1}^{T}\boldsymbol{Z}^{(l)} = \begin{bmatrix} (1+\eta\tau)^{l}\boldsymbol{U}_{1}\boldsymbol{A}_{1} & \boldsymbol{U}_{1}\boldsymbol{E}_{2,1} & \dots & \boldsymbol{U}_{1}\boldsymbol{E}_{K,1} \end{bmatrix}, \\ \boldsymbol{M}_{1}^{(l)} = \begin{bmatrix} (1+\eta\tau)^{2l}\boldsymbol{A}_{1}^{T}\boldsymbol{A}_{1} & (1+\eta\tau)^{l}\boldsymbol{A}_{1}^{T}\boldsymbol{E}_{2,1} & \dots & (1+\eta\tau)^{l}\boldsymbol{A}_{1}^{T}\boldsymbol{E}_{K,1} \\ (1+\eta\tau)^{l}\boldsymbol{E}_{2,1}^{T}\boldsymbol{A}_{1} & \boldsymbol{E}_{2,1}^{T}\boldsymbol{E}_{2,1} & \dots & \boldsymbol{E}_{2,1}^{T}\boldsymbol{E}_{K,1} \\ \vdots & \vdots & \ddots & \vdots \\ (1+\eta\tau)^{l}\boldsymbol{E}_{K,1}^{T}\boldsymbol{A}_{1} & \boldsymbol{E}_{K,1}^{T}\boldsymbol{E}_{2,1} & \dots & \boldsymbol{E}_{K,1}^{T}\boldsymbol{E}_{K,1} \end{bmatrix}$$

Using the same argument, we can compute $M_k^{(l)}$ for each $k \in [K]$. This, together with (30) for each $k \in [K]$, yields

1023
1024
$$\sum_{k=1}^{K} U_k U_k^T Z^{(0)} \varphi(M_k^{(0)}) = \left[(1+\eta\tau)^l \tau U_1 A_1 \quad (1+\eta\tau)^l \tau U_2 A_2 \quad \dots \quad (1+\eta\tau)^l \tau U_K A_K \right].$$
1025

Using this, (27), and (3), we directly obtain that (35) holds for l + 1. Then, we prove the claim.

1026 1.2 1.2Llama Least Squares 1027 1.0 Llama AoT 1.0 Llama 1028 Squared Error Squared Error Least Squares Llama AoT 0.8 1029 0.8 Lasso (alpha=0.01) 1030 0.6 0.6 1031 0.40.41032 0.2 0.2 1033 1034 0.0 0.0 1035 0 10 20 30 40 0 10 20 30 40 1036 In-context Examples In-context Examples 1037 (a) Linear regression (b) Sparse linear regression

Figure 7: Evaluating models of Llama architectures on in-context learning linear functions. We plot the normalized squared error as a function of in-context examples.

1042 B SUPPLEMENTARY EXPERIMENTS

1044 B.0.1 MORE ON ICL

In addition, we performed the same ICL analysis as in Section 4.1.2. All the settings are the same, except that we changed the base model architecture to Llama (Touvron et al., 2023). And, we can see that the results are similar.

1049

1039

1040 1041

1043

1045

B.0.2 Emergence of Semantic Properties

The attention heads in our models have different semantic meanings, and indeed demonstrate the interpretability of our proposed architecture in practice. In Figure 8, we visualize the self-attention heatmaps between the [CLS] token and other image patches. We select 5 attention heads by manual inspection and find that they capture different parts of objects, displaying different semantic meanings.

1056 1057 B.O.3 COMPUTING REQUIREMENT

In this section, we present the computing resources of a forward pass used by AoT-based language models and GPT-2 empirically in Table 3. The context window is 1024 tokens and the batch size is 16. The GFLOPS is measured by the PyTorch profiler, the total GPU memory consumption by the NVIDIA System Management Interface, and the running time of one forward pass by the Python time module. The only optimization we use is the default mode of the PyTorch compiler.

Table 3: The GFLOPS, total GPU memory consumption, and the running time of one forward pass are shown of AoT and GPT-2 at different sizes.

1066	Models	GFLOPS	Total GPU Memory in MiB	Running time in ms
1062	Base 102M	1651	21482	43
1069	Medium 182M	3868	36198	78
1070	Large 326M	8056	57896	225
1071	GPT-2 Base 124M	2785	23300	27
1072	GPT-2 Medium 335M	9898	51578	158

1073 1074

1063

1075

1076

1077 1078



Figure 8: **Visualization of attention heads.** We feed our AoT a mini-batch of images and extract the attention maps of different heads from the penultimate layer. We show that these heads capture certain semantic meanings across different images.