

On the Optimization and Generalization of Multi-head Attention

Anonymous authors

Paper under double-blind review

Abstract

The training and generalization dynamics of the Transformer’s core mechanism, namely the Attention mechanism, remain under-explored. Besides, existing analyses primarily focus on single-head attention. Inspired by the demonstrated benefits of overparameterization when training fully-connected networks, we investigate the potential optimization and generalization advantages of using multiple attention heads. Towards this goal, we derive convergence and generalization guarantees for gradient-descent training of a single-layer multi-head self-attention model, under a suitable realizability condition on the data. We then establish primitive conditions on the initialization that ensure realizability holds. Finally, we demonstrate that these conditions are satisfied for a simple tokenized-mixture model. We expect the analysis can be extended to various data-model and architecture variations.

1 Introduction

Transformers have emerged as a promising paradigm in deep learning, primarily attributable to their distinctive self-attention mechanism. Motivated by the model’s state-of-the-art performance in natural language processing (Devlin et al., 2019; Brown et al., 2020; Raffel et al., 2020) and computer vision (Dosovitskiy et al., 2021; Radford et al., 2021; Touvron et al., 2021), the theoretical study of the attention mechanism has seen a notable surge in interest recently. Numerous studies have already explored the expressivity of Attention, e.g. (Baldi & Vershynin, 2022; Dong et al., 2021; Yun et al., 2020a,b; Sanford et al., 2023; Bietti et al., 2023), and initial findings regarding memory capacity have been very recently studied in (Baldi & Vershynin, 2022; Dong et al., 2021; Yun et al., 2020a,b; Mahdavi et al., 2023). In an attempt to comprehend optimization aspects of training attention models, Sahiner et al. (2022); Ergen et al. (2022) have investigated convex-relaxations, while Tarzanagh et al. (2023a) investigates the model’s implicit bias. Additionally, Edelman et al. (2021) have presented capacity and Rademacher complexity-based generalization bounds for Self-Attention. However, the exploration of the *finite-time* optimization and generalization dynamics of gradient-descent (GD) for training attention models largely remains an open question.

Recent contributions in this direction, which serve as motivation for our work, include the studies by (Jelassi et al., 2022; Li et al., 2023a; Oymak et al., 2023). These works concentrate on single-layer attention models with a *single attention head*. Furthermore, despite necessary simplifying assumptions made for the data, the analyses in all three cases are rather intricate and appear highly specialized on the individual attention and data model. These direct and highly specialized analyses present certain challenges. First, it remains uncertain whether they can be encompassed within a broader framework that can potentially be extended to more complex attention architectures and diverse data models. Second, they appear disconnected from existing frameworks that have been flourishing in recent years for conventional architectures like fully-connected and convolutional neural networks e.g., (Jacot et al., 2018; Ji & Telgarsky, 2020; Richards & Rabbat, 2021; Banerjee et al., 2022; Taheri & Thrampoulidis, 2023). Consequently, it is also unclear how the introduction of attention alters the analysis landscape.

In this work, we study the optimization and generalization properties of multi-head attention mechanism trained by gradient methods. Our approach specifically leverages the use of *multiple attention heads*. Despite the operational differences between attention heads in an attention model and hidden nodes in an MLP,

we demonstrate, from an analysis perspective, that this parallelism enables the exploitation of frameworks developed for the latter to study the former. Particularly for the generalization analysis, we leverage recent advancements in the application of the algorithmic-stability framework to overparameterized MLPs (Richards & Kuzborskij, 2021; Taheri & Thrampoulidis, 2023).

Contributions. We study training and generalization of gradient descent optimization for a multi-head attention (MHA) layer with H heads in a binary classification setting. For this setting, detailed in Section 2, we analyze training with logistic loss both the attention weights (parameterizing the softmax logits), as well as, the linear decoder that turns output tokens to label prediction.

In Section 3, we characterize key properties of the empirical loss \widehat{L} , specifically establishing that it is self-bounded and satisfies a key self-bounded weak-convexity property, i.e. $\lambda_{\min}(\nabla^2 \widehat{L}(\boldsymbol{\theta})) \gtrsim -\frac{\kappa}{\sqrt{H}} \widehat{L}(\boldsymbol{\theta})$ for a parameter κ that depends only mildly on the parameter vector $\boldsymbol{\theta}$. Establishing these properties (and also quantifying κ) involves carefully computing and bounding the gradient and Hessian of the MHA layer, calculations that can be useful beyond the context of our paper.

In Sections 4.1-4.2, we present our training and generalization bounds in their most general form. The bounds are given in terms of the empirical loss $\widehat{L}(\boldsymbol{\theta})$ and the distance $\|\boldsymbol{\theta} - \boldsymbol{\theta}_0\|$ to initialization $\boldsymbol{\theta}_0$ of an appropriately chosen target vector $\boldsymbol{\theta}$. The distance to initialization also controls the minimum number of heads $H \gtrsim \|\boldsymbol{\theta} - \boldsymbol{\theta}_0\|^6$ required for the bounds to hold. The choice of an appropriate parameter $\boldsymbol{\theta}$ that makes the bounds tight is generically specific to the data setting and the chosen initialization. To guide such a choice, in Section 4.3, we formalize primitive and straightforward-to-check conditions on the initialization $\boldsymbol{\theta}_0$ that ensure it is possible to find an appropriate $\boldsymbol{\theta}$. In short, provided the model output at initialization is logarithmic on the train-set size n and the data are separable with respect to the neural-tangent kernel (NTK) features of the MHA model with constant margin γ , then Corollary 2 shows that with step-size $\eta = \widetilde{O}(1)$ and $\Theta(n)$ gradient descent steps, the train loss and generalization gap is bounded by $\widetilde{O}(1/n)$ provided only a polylogarithmic number of heads $H = \Omega(\log^6(n))$. We remark that the aforementioned NTK separability assumption, although related to, differs from the standard NTK analysis. Besides, while this assumption is sufficient to apply our general bounds, it is not a necessary condition.

In Section 5, we investigate a tokenized mixture data model with label-(ir)relevant tokens. We show that after one randomized gradient step from zero initialization, the NTK features of the MHA model separate the data with margin γ_* . Thus, applying our general analysis from Section 4.1 we establish training and generalization bounds as described above, for a logarithmic number of heads. Towards assessing the optimality of these bounds, we demonstrate that MHA is expressive enough to achieve margin γ_{attn} that is superior to γ_* . The mechanism to reach γ_{attn} involves selecting key-query weights of sufficiently large norm, which saturates the softmax nonlinearity by suppressing label-irrelevant tokens. We identify the large-norm requirement as a potential bottleneck in selecting those weights as target parameters in our theory framework and discuss open questions regarding extending the analytical framework into this specific regime.

The remaining parts are organised as follows. Proof sketches of our main training/generalization bounds are given in Section 6. The paper concludes in Section 7 with remarks on our findings' implications and open questions. Detailed proofs are in the appendix, where we also present synthetic numerical experiments.

Related work. We give a brief overview of the most relevant works on understanding optimization/generalization of self-Attention or its variants. Please see Section H for more detailed exposition. Oymak et al. (2023) diverges from traditional self-Attention by focusing on a variant called prompt-Attention, aiming to gain understanding of prompt-tuning. Jelassi et al. (2022) shed light on how ViTs learn spatially localized patterns using gradient-based methods. Li et al. (2023a) provides sample complexity bounds for achieving zero generalization error on training three-layer ViTs for classification tasks for a similar tokenized mixture data model as ours. Contemporaneous work Tian et al. (2023) presents SGD-dynamics of single-layer attention for next-token prediction by re-parameterizing the original problem in terms of the softmax and classification logit matrices, while Tarzanagh et al. (2023b,a) study the implicit bias of training the softmax weights \mathbf{W} with a fixed decoder \mathbf{U} . All these works focus on a single attention head; instead, we leverage the use of multiple heads to establish connections to the literature on GD training of overparameterized MLPs. Conceptually, Hron et al. (2020) drew similar connections, linking multi-head attention to a Gaussian process

in the limit as the number of heads approaches infinity. In contrast, we study the more practical regime of finite heads and obtain *finite-time* optimization and generalization bounds.

Among the extensive studies on training/generalization of overparameterized MLPs, our work closely aligns with [Nitanda et al. (2019); Ji & Telgarsky (2020); Cao & Gu (2019); Chen et al. (2020); Telgarsky (2022); Taheri & Thrampoulidis (2023)] focusing on classification with logistic loss. Conceptually, our findings extend this research to attention models. The use of algorithmic-stability tools towards order-optimal generalization bounds for overparameterized MLPs has been exploited recently by [Richards & Kuzborskij (2021); Richards & Rabbat (2021); Taheri & Thrampoulidis (2023); Lei et al. (2022)]. To adapt these tools to the MHA layer, we critically utilize the smoothness of the softmax function and derive bounds on the growth of the model’s gradient/Hessian, which establish a self-bounded weak convexity property for the empirical risk (see Corollary 1). Our approach also involves training both the classifier and attention weights, necessitating several technical adjustments detailed in Section 6 and Appendix B.1

2 Preliminaries

Notation. $\varphi(\cdot) : \mathbb{R}^T \rightarrow \mathbb{R}^T$ denotes the softmax map and $\varphi'(\mathbf{v}) := \nabla \varphi(\mathbf{v}) = \text{diag}(\varphi(\mathbf{v})) - \mathbf{v}\mathbf{v}^\top$ its gradient at $\mathbf{v} \in \mathbb{R}^T$. For $t \in [T]$, $\varphi_t(\mathbf{v})$ is the t -th entry of $\varphi(\mathbf{v}) \in \mathbb{R}^T$. For $\mathbf{A} \in \mathbb{R}^{n \times m}$, $\mathbf{A}_{i,:}$ is its i -th row and $\mathbf{A}_{:,j}$ is its j -th column. Recall the induced matrix norm $\|\mathbf{A}\|_{p,q} = \max_{\|\mathbf{v}\|_p=1} \|\mathbf{A}\mathbf{v}\|_q$ and particularly the following: $\|\mathbf{A}\|_{2,\infty} = \max_{i \in [n]} \|\mathbf{A}_{i,:}\|$, $\|\mathbf{A}\|_{1,2} = \max_{j \in [m]} \|\mathbf{A}_{:,j}\|$, and $\|\mathbf{A}\|_{1,\infty} = \max_{j \in [m]} \|\mathbf{A}_{:,j}\|_\infty$. For simplicity, $\|\mathbf{A}\|$, $\|\mathbf{v}\|$ denote Euclidean norms and $\lambda_{\min}(\mathbf{A})$ the minimum eigenvalue. We let $a \wedge b = \min\{a, b\}$ and $a \vee b = \max\{a, b\}$. concat denotes vector concatenation. All logarithms are natural logarithms (base e). We represent the line segment between $\mathbf{w}_1, \mathbf{w}_2 \in \mathbb{R}^{d'}$ as $[\mathbf{w}_1, \mathbf{w}_2] = \{\mathbf{w} : \mathbf{w} = \alpha \mathbf{w}_1 + (1 - \alpha) \mathbf{w}_2, \alpha \in [0, 1]\}$. Finally, to simplify the exposition we use “ \gtrsim ” or “ \lesssim ” notation to hide absolute constants. We also occasionally use standard notations \mathcal{O}, Ω and $\tilde{\mathcal{O}}, \tilde{\Omega}$ to hide poly-logarithmic factors. Unless otherwise stated these order-wise notations are with respect to the training-set size n . Whenever used, exact constants are specified in the appendix.

Single-head Self-attention. A single-layer self-attention head $\text{ATTN} : \mathbb{R}^{T \times d} \rightarrow \mathbb{R}^{T \times d}$ with context size T and dimension d parameterized by key, query and value matrices $\mathbf{W}_Q, \mathbf{W}_K \in \mathbb{R}^{d \times d_h}$, $\mathbf{W}_V \in \mathbb{R}^{d \times d_v}$ is given by:

$$\text{ATTN}(\mathbf{X}; \mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V) := \varphi(\mathbf{X} \mathbf{W}_Q \mathbf{W}_K^\top \mathbf{X}^\top) \mathbf{X} \mathbf{W}_V.$$

Here, $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T]^\top \in \mathbb{R}^{T \times d}$ is the input token matrix and $\varphi(\mathbf{X} \mathbf{W}_Q \mathbf{W}_K^\top \mathbf{X}^\top) \in \mathbb{R}^{T \times T}$ is the attention matrix. (Softmax applied to a matrix acts row-wise.) To turn the Attention output in a prediction label, we compose ATTN with a linear projection head (aka decoder). Thus, the model’s output is*

$$\Phi(\mathbf{X}; \mathbf{W}, \mathbf{U}) := \langle \mathbf{U}, \varphi(\mathbf{X} \mathbf{W} \mathbf{X}^\top) \mathbf{X} \rangle. \quad (1)$$

Note that we absorb the value weight matrix \mathbf{W}_V into the projector $\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_T]^\top \in \mathbb{R}^{T \times d}$. Also, we parameterize throughout the key-query product matrix as $\mathbf{W} := \mathbf{W}_Q \mathbf{W}_K^\top$.

Multi-head Self-attention. Our focus is on the multi-head attention (MHA) model with H heads:

$$\sum_{h \in [H]} \text{ATTN}(\mathbf{X}; \mathbf{W}_{Q_h}, \mathbf{W}_{K_h}, \mathbf{W}_{V_h}) \mathbf{W}_{O_h},$$

for output matrices $\mathbf{W}_{O_h} \in \mathbb{R}^{d_v \times d}$. Absorbing $\mathbf{W}_{V_h} \mathbf{W}_{O_h}$ into a projection layer (similar to the single-head attention) and parameterizing $\mathbf{W}_h := \mathbf{W}_{Q_h} \mathbf{W}_{K_h}^\top$ we arrive at the following MHA model:

$$\tilde{\Phi}(\mathbf{X}; \tilde{\mathbf{W}}, \tilde{\mathbf{U}}) := \frac{1}{\sqrt{H}} \sum_{h \in [H]} \Phi(\mathbf{X}; \mathbf{W}_h, \mathbf{U}_h) = \frac{1}{\sqrt{H}} \sum_{h \in [H]} \langle \mathbf{U}_h, \varphi(\mathbf{X} \mathbf{W}_h \mathbf{X}^\top) \mathbf{X} \rangle, \quad (2)$$

parameterized by $\tilde{\mathbf{W}} := \text{concat}(\{\mathbf{W}_h\}_{h \in [H]})$ and $\tilde{\mathbf{U}} := \text{concat}(\{\mathbf{U}_h\}_{h \in [H]})$. The $1/\sqrt{H}$ scaling is analogous to the normalization in MLP literature e.g. [Du et al. 2019; Ji & Telgarsky, 2021; Richards & Kuzborskij, 2021], ensuring the model variance is of constant order when \mathbf{U}_h is initialized $\mathcal{O}_H(1)$.

*While we focus on (i) Full-projection: trainable matrix $\mathbf{U} \in \mathbb{R}^{T \times d}$, our results also apply to (ii) Pooling: $\mathbf{U} = \mathbf{u} \mathbf{1}_T^\top$ with trainable $\mathbf{u} \in \mathbb{R}^d$, and (iii) Last-token output: $\mathbf{U} = [0_{d \times (T-1)} \quad \mathbf{u}]^\top$ with trainable $\mathbf{u} \in \mathbb{R}^d$.

†These relaxations sacrifice some generality since it is common practice to set values for d_h and d_v such that $d_v = d/H < d$, thus imposing low-rank restrictions on the product matrices $\mathbf{W}_{Q_h} \mathbf{W}_{K_h}^\top$, $\mathbf{W}_{V_h} \mathbf{W}_{O_h}$. We defer a treatment of these to future work.

Throughout, we will use $\boldsymbol{\theta}_h := \text{concat}(\mathbf{U}_h, \mathbf{W}_h) \in \mathbb{R}^{dT+d^2}$, to denote the trainable parameters of the h -attention head and $\tilde{\boldsymbol{\theta}} := \text{concat}(\{\boldsymbol{\theta}_h\}_{h \in [H]}) \in \mathbb{R}^{H(dT+d^2)}$ for the trainable parameters of the overall model. More generally, we use the convention of applying “ \sim ” notation for quantities relating to the multi-head model. Finally, with some slight abuse of notation, we define: $\|\tilde{\boldsymbol{\theta}}\|_{2,\infty} := \max_{h \in [H]} \|\boldsymbol{\theta}_h\|$.

Training. Given training set $(\mathbf{X}_i, y_i)_{i \in [n]}$, with n IID samples, we minimize logistic-loss based empirical risk

$$\widehat{L}(\tilde{\boldsymbol{\theta}}) := \frac{1}{n} \sum_{i \in [n]} \ell(y_i \tilde{\Phi}(\mathbf{X}_i; \tilde{\boldsymbol{\theta}})) := \frac{1}{n} \sum_{i \in [n]} \log(1 + e^{-y_i \tilde{\Phi}(\mathbf{X}_i; \tilde{\boldsymbol{\theta}})}).$$

Our analysis extends to any convex, smooth, Lipschitz and self-bounded loss.[‡] The empirical risk is minimized as an approximation of the *test loss* defined as $L(\tilde{\boldsymbol{\theta}}) := \mathbb{E}_{(\mathbf{X}, y)}[\ell(y \tilde{\Phi}(\mathbf{X}; \tilde{\boldsymbol{\theta}}))]$. We consider standard gradient-descent (GD) applied to empirical risk \widehat{L} . Formally, initialized at $\tilde{\boldsymbol{\theta}}^{(0)}$ and equipped with step-size $\eta > 0$, at each iteration $k \geq 0$, GD performs the following update:

$$\tilde{\boldsymbol{\theta}}^{(k+1)} = \tilde{\boldsymbol{\theta}}^{(k)} - \eta \nabla \widehat{L}(\tilde{\boldsymbol{\theta}}^{(k)}).$$

3 Gradient and Hessian bounds of soft-max attention

This section establishes bounds on the gradient and Hessian of the logistic empirical risk $\widehat{L}(\cdot)$ evaluated on the multi-head attention model. To do this, we first derive bounds on the Euclidean norm and spectral-norm for the gradient and Hessian of the self-attention model. In order to simplify notations, we state here the bounds for the single-head model (see App. A.1 for multi-head model): $\Phi(\mathbf{X}; \boldsymbol{\theta}) := \Phi(\mathbf{X}; \mathbf{W}, \mathbf{U}) = \langle \mathbf{U}, \varphi(\mathbf{X} \mathbf{W} \mathbf{X}^\top) \mathbf{X} \rangle$.

Lemma 1 (Gradient/Hessian formulas). *For all $\mathbf{a} \in \mathbb{R}^T$, $\mathbf{b}, \mathbf{c} \in \mathbb{R}^d$ the model’s gradients satisfy:*

$$\bullet \quad \nabla_{\mathbf{U}} \Phi(\mathbf{X}; \boldsymbol{\theta}) = \varphi(\mathbf{X} \mathbf{W} \mathbf{X}^\top) \mathbf{X}, \quad \text{and} \quad \nabla_{\mathbf{W}} \Phi(\mathbf{X}; \boldsymbol{\theta}) = \sum_{t=1}^T \mathbf{x}_t \mathbf{u}_t^\top \mathbf{X}^\top \varphi'(\mathbf{X} \mathbf{W}^\top \mathbf{x}_t) \mathbf{X}.$$

$$\bullet \quad \nabla_{\mathbf{W}} \langle \mathbf{a}, \nabla_{\mathbf{U}} \Phi(\mathbf{X}; \boldsymbol{\theta}) \mathbf{b} \rangle = \sum_{t=1}^T \mathbf{x}_t \mathbf{a}_t^\top \mathbf{b}^\top \mathbf{X}^\top \varphi'(\mathbf{X} \mathbf{W}^\top \mathbf{x}_t) \mathbf{X}, \quad \text{and}$$

$$\nabla_{\mathbf{W}} \langle \mathbf{c}, \nabla_{\mathbf{W}} \Phi(\mathbf{X}; \boldsymbol{\theta}) \mathbf{b} \rangle = \sum_{t=1}^T (\mathbf{c}^\top \mathbf{x}_t) \mathbf{x}_t \mathbf{d}^\top \varphi'(\mathbf{X} \mathbf{W}^\top \mathbf{x}_t) \mathbf{X}$$

$$\text{where } \mathbf{d} := \text{diag}(\mathbf{X} \mathbf{b}) \mathbf{X} \mathbf{u}_t - \mathbf{X} \mathbf{u}_t \mathbf{b}^\top \mathbf{X}^\top \varphi(\mathbf{X} \mathbf{W}^\top \mathbf{x}_t) - \mathbf{X} \mathbf{b} \mathbf{u}_t^\top \mathbf{X}^\top \varphi(\mathbf{X} \mathbf{W}^\top \mathbf{x}_t).$$

These calculations imply the following useful bounds.

Proposition 1 (Model Gradient/Hessian bounds). *The Euclidean norm of the gradient and the spectral norm of the Hessian of the single-head Attention model (I) are bounded as follows:*

$$\bullet \quad \|\nabla_{\boldsymbol{\theta}} \Phi(\mathbf{X}; \boldsymbol{\theta})\| \leq 2 \|\mathbf{X}\|_{2,\infty}^2 \sum_{t=1}^T \|\mathbf{X} \mathbf{u}_t\|_\infty + \sqrt{T} \|\mathbf{X}\|_{2,\infty}.$$

$$\bullet \quad \|\nabla_{\boldsymbol{\theta}}^2 \Phi(\mathbf{X}; \boldsymbol{\theta})\| \leq 6d \|\mathbf{X}\|_{2,\infty}^2 \|\mathbf{X}\|_{1,\infty}^2 \sum_{t=1}^T \|\mathbf{X} \mathbf{u}_t\|_\infty + 2\sqrt{Td} \|\mathbf{X}\|_{2,\infty}^2 \|\mathbf{X}\|_{1,\infty}.$$

Next, we focus on the empirical loss \widehat{L} . To derive bounds on its gradient and Hessian, we leverage the model’s bounds from Proposition 1 and the fact that logistic loss is self-bounded, i.e., $|\ell'(t)| \leq \ell(t)$. To provide concrete statements, we introduce first a mild boundedness assumption.

Assumption 1 (Bounded data). *Data $(\mathbf{X}, y) \in \mathbb{R}^{T \times d} \times \mathbb{R}$ satisfy the following conditions almost surely: $y \in \{\pm 1\}$, and for some $R \geq 1$, it holds for all $t \in [T]$ that $\|\mathbf{x}_t\| \leq R$.*

Corollary 1 (Loss properties). *Under Assumption 1, the objective’s gradient and Hessian satisfy the bounds.[§]*

[‡]A function $\ell : \mathbb{R} \rightarrow \mathbb{R}$ is self-bounded if $\exists C > 0$ such that $|\ell'(t)| \leq C \ell(t)$.

[§]In all the bounds in this paper involving $\|\tilde{\boldsymbol{\theta}}\|_{2,\infty}$, it is possible to substitute this term with $\max_{h \in [H]} \|\mathbf{U}_h\|$. However, for the sake of notation simplicity, we opt for a slightly looser bound $\max_{h \in [H]} \|\mathbf{U}_h\| \leq \max_{h \in [H]} \|\boldsymbol{\theta}_h\| =: \|\tilde{\boldsymbol{\theta}}\|_{2,\infty}$.

$$\begin{aligned}
(1) \quad & \|\nabla \widehat{L}(\tilde{\boldsymbol{\theta}})\| \leq \beta_1(\tilde{\boldsymbol{\theta}}) \widehat{L}(\tilde{\boldsymbol{\theta}}), & \beta_1(\tilde{\boldsymbol{\theta}}) &:= \sqrt{T} R (2 R^2 \|\tilde{\boldsymbol{\theta}}\|_{2,\infty} + 1). \\
(2) \quad & \|\nabla^2 \widehat{L}(\tilde{\boldsymbol{\theta}})\| \leq \beta_2(\tilde{\boldsymbol{\theta}}), & \beta_2(\tilde{\boldsymbol{\theta}}) &:= \frac{1}{\sqrt{H}} \beta_3(\tilde{\boldsymbol{\theta}}) + \frac{1}{4} \beta_1(\tilde{\boldsymbol{\theta}})^2. \\
(3) \quad & \lambda_{\min}(\nabla^2 \widehat{L}(\tilde{\boldsymbol{\theta}})) \geq -\frac{\beta_3(\tilde{\boldsymbol{\theta}})}{\sqrt{H}} \widehat{L}(\tilde{\boldsymbol{\theta}}) & \beta_3(\tilde{\boldsymbol{\theta}}) &:= 2 \sqrt{T} d R^3 (3 \sqrt{d} R^2 \|\tilde{\boldsymbol{\theta}}\|_{2,\infty} + 1).
\end{aligned}$$

The loss properties above are crucial for the training and generalization analysis. Property (1) establishes self-boundedness of the empirical loss, which is used to analyze stability of GD updates for generalization. Property (2) is used to establish descent of gradient updates for appropriate choice of step-size η . Note that the smoothness upper bound is $\tilde{\boldsymbol{\theta}}$ -dependent, hence to show descent we need to also guarantee boundedness of the updates. Finally, property (3) establishes a self-bounded weak-convexity property of the loss, which is crucial to both the training and generalization analysis. Specifically, as the number of heads H increases, the minimum eigenvalue becomes less positive, indicating an approach towards convex-like behavior.

4 Main results

In this section, we present our training and generalization bounds for multi-head attention.

4.1 Training bounds

We state our main result on train loss convergence in the following theorem. See App. B for exact constants and the detailed proofs.

Theorem 1 (Training loss). *Fix iteration horizon $K \geq 1$ and any $\tilde{\boldsymbol{\theta}} \in \mathbb{R}^{H(dT+d^2)}$ and H satisfying*

$$\sqrt{H} \gtrsim d T^{1/2} R^5 \|\tilde{\boldsymbol{\theta}}\|_{2,\infty} \|\tilde{\boldsymbol{\theta}} - \tilde{\boldsymbol{\theta}}^{(0)}\|^3. \quad (3)$$

Fix step-size $\eta \leq 1 \wedge 1/\rho(\tilde{\boldsymbol{\theta}}) \wedge \frac{\|\tilde{\boldsymbol{\theta}} - \tilde{\boldsymbol{\theta}}^{(0)}\|^2}{K \widehat{L}(\tilde{\boldsymbol{\theta}})} \wedge \frac{\|\tilde{\boldsymbol{\theta}} - \tilde{\boldsymbol{\theta}}^{(0)}\|^2}{\widehat{L}(\tilde{\boldsymbol{\theta}}^{(0)})}$, with $\rho(\tilde{\boldsymbol{\theta}}) \lesssim d^{3/2} T^{3/2} R^{13} \|\tilde{\boldsymbol{\theta}}\|_{2,\infty}^2 \|\tilde{\boldsymbol{\theta}} - \tilde{\boldsymbol{\theta}}^{(0)}\|^2$. Then, the following bounds hold for the training loss and the weights' norm at iteration K of GD:

$$\begin{aligned}
\widehat{L}(\tilde{\boldsymbol{\theta}}^{(K)}) &\leq \frac{1}{K} \sum_{k=1}^K \widehat{L}(\tilde{\boldsymbol{\theta}}_k) \leq 2 \widehat{L}(\tilde{\boldsymbol{\theta}}) + \frac{5 \|\tilde{\boldsymbol{\theta}} - \tilde{\boldsymbol{\theta}}^{(0)}\|^2}{4 \eta K}, \\
\|\tilde{\boldsymbol{\theta}}^{(K)} - \tilde{\boldsymbol{\theta}}^{(0)}\| &\leq 4 \|\tilde{\boldsymbol{\theta}} - \tilde{\boldsymbol{\theta}}^{(0)}\|.
\end{aligned} \quad (4)$$

Yielding a concrete train loss bound requires an appropriate set of target parameters $\tilde{\boldsymbol{\theta}}$ in the sense of minimizing the bound in (4). Hence, $\tilde{\boldsymbol{\theta}}$ should simultaneously attain small loss ($\widehat{L}(\tilde{\boldsymbol{\theta}})$) and distance to initialization ($\|\tilde{\boldsymbol{\theta}} - \tilde{\boldsymbol{\theta}}^{(0)}\|$). This desiderata is formalized in Assumption 2 below. The distance to initialization, as well as $\|\tilde{\boldsymbol{\theta}}\|_{2,\infty}$, determine how many heads are required for our bounds to hold. Also, in view of the bound in (4), it is reasonable that an appropriate choice for $\tilde{\boldsymbol{\theta}}$ attains $\widehat{L}(\tilde{\boldsymbol{\theta}})$ of same order as $\|\tilde{\boldsymbol{\theta}} - \tilde{\boldsymbol{\theta}}^{(0)}\|^2/K$. Hence, the theorem's restriction on the step-size is governed by the inverse local-smoothness of the loss: $\eta \lesssim 1/\rho(\tilde{\boldsymbol{\theta}})$.

4.2 Generalization bounds

Next we bound the expected generalization gap. Expectations are with respect to (w.r.t) randomness of the train set. See App. C for the detailed proof, which is based on algorithmic-stability.

Theorem 2 (Generalization loss). *Fix any $K \geq 1$, any $\tilde{\boldsymbol{\theta}}$ and H satisfying (3), and any step-size η satisfying the conditions of Thm. 1. Then the expected generalization gap at iteration K satisfies,*

$$\mathbb{E}[L(\tilde{\boldsymbol{\theta}}^{(K)}) - \widehat{L}(\tilde{\boldsymbol{\theta}}^{(K)})] \leq \frac{4}{n} \mathbb{E}\left[2 K \widehat{L}(\tilde{\boldsymbol{\theta}}) + \frac{9 \|\tilde{\boldsymbol{\theta}} - \tilde{\boldsymbol{\theta}}^{(0)}\|^2}{4 \eta}\right]. \quad (5)$$

The condition on the number of heads is same up to constants to the corresponding condition in Theorem 1. Also, the generalization-gap bound translates to test-loss bound by combining with Thm. 1. Finally, similar

to Thm. [1](#), we can get concrete bounds under the realizability assumption; see Cor. [4](#) in App. [C.2](#). For the generalization analysis, we require that the realizability assumption holds almost surely over all training sets sampled from the data distribution.

The bounds on optimization and generalization are up to constants same as analogous bounds for logistic regression ([Soudry et al., 2018](#); [Ji & Telgarsky, 2018](#); [Shamir, 2021](#)). Yet, for these bounds to be valid, we require sufficiently large number of heads as well as the existence of an appropriate set of target parameters $\tilde{\theta}$, as stated in the conditions of theorem. Namely, these conditions are related to the realizability condition, which guarantees small training error near initialization. The next assumption formalizes these conditions.

Assumption 2 (Realizability). *There exist non-increasing functions $g : \mathbb{R}_+ \rightarrow \mathbb{R}_+$ and $g_0 : \mathbb{R}_+ \rightarrow \mathbb{R}_+$ such that $\forall \epsilon > 0$, there exists model parameters $\tilde{\theta}_{(\epsilon)} \in \mathbb{R}^{H(dT+d^2)}$ for which: (i) the empirical loss over n data samples satisfies $\widehat{L}(\tilde{\theta}_{(\epsilon)}) \leq \epsilon$, (ii) $\|\tilde{\theta}_{(\epsilon)} - \tilde{\theta}^{(0)}\| \leq g_0(\epsilon)$, and, (iii) $\|\tilde{\theta}_{(\epsilon)}\|_{2,\infty} \leq g(\epsilon)$.*

With this assumption, we can specialize the result of Thms. above to specific data settings; see Cor. [3](#) and [4](#) in App. [B.5](#) and [C.2](#). In the next section we will further show how the realizability assumption is satisfied.

4.3 Primitive conditions for checking realizability

Here, we introduce a set of more primitive and straightforward-to-check conditions on the data and initialization that ensure the realizability Assumption [2](#) holds.

Definition 1 (Good initialization). *We say $\tilde{\theta}^{(0)} = \text{concat}(\theta_1^{(0)}, \dots, \theta_H^{(0)})$ is a **good initialization** with respect to training data $(\mathbf{X}_i, y_i)_{i \in [n]}$ provided the following three properties hold.*

- P1. **Parameter $L_{2,\infty}$ -bound:** There exists parameter $B_2 \geq 1$ such that $\forall h \in [H]$ it holds $\|\theta_h^{(0)}\|_2 \leq B_2$.*
- P2. **Model bound:** There exists parameter $B_\Phi \geq 1$ such that $\forall i \in [n]$ it holds $|\tilde{\Phi}(\mathbf{X}_i; \tilde{\theta}^{(0)})| \leq B_\Phi$.*
- P3. **NTK separability:** There exists $\tilde{\theta}_* \in \mathbb{R}^{H(dT+d^2)}$ and $\gamma > 0$ such that $\|\tilde{\theta}_*\| = \sqrt{2}$ and $\forall i \in [n]$, it holds $y_i \langle \nabla \tilde{\Phi}(\mathbf{X}_i; \tilde{\theta}^{(0)}), \tilde{\theta}_* \rangle \geq \gamma$.*

Prop. [7](#) in the appendix shows that starting from a **good initialization** we can always find $\tilde{\theta}_{(\epsilon)}$ satisfying the realizability Assumption [2](#) provided large enough number of heads. Thus, given **good initialization**, we can immediately apply Theorems [1](#) and [2](#) to get the following concrete bounds.

Corollary 2 (General bounds under good initialization). *Suppose **good initialization** $\tilde{\theta}^{(0)}$ and let*

$$\sqrt{H} \gtrsim dT^{1/2} R^5 B_2^2 (g_0(1/K))^3, \quad \text{where } g_0\left(\frac{1}{K}\right) = \frac{2B_\Phi + \log(K)}{\gamma}.$$

Further fix step-size $\eta \leq 1 \wedge 1/\rho(K) \wedge \frac{4B_\Phi^2}{\gamma^2 \log(1+e^{B_\Phi})}$ with $\rho(K) \gtrsim d^{3/2} T^{3/2} R^{13} g_0(\frac{1}{K})^4$. Then, it holds that

$$\widehat{L}(\tilde{\theta}^{(K)}) \leq \frac{2}{K} + \frac{5(2B_\Phi + \log(K))^2}{4\gamma^2 \eta K}, \quad \text{and} \quad \mathbb{E}[L(\tilde{\theta}^{(K)}) - \widehat{L}(\tilde{\theta}^{(K)})] \leq \frac{17(2B_\Phi + \log(K))^2}{\gamma^2 \eta n}.$$

Consider training loss after K GD steps: Assuming $B_\Phi = \tilde{\mathcal{O}}_K(1)$ and $\gamma = \mathcal{O}_K(1)$, then choosing $\eta = \tilde{\mathcal{O}}_K(1)$, the corollary guarantees train loss is $\tilde{\mathcal{O}}_K(\frac{1}{K})$ provided polylogarithmic number of heads $H = \Omega(\log^6(K))$. Moreover, after $K \approx n$ GD steps the expected test loss is $\mathcal{O}(1/n)$.

Remark 1. *The last two conditions (P2 and P3) for **good initialization** are similar to the conditions needed in ([Taheri & Thrampoulidis, 2023](#); [Ji & Telgarsky, 2020](#); [Nitanda et al., 2019](#)) for analysis of two-layer MLPs. Compared to ([Ji & Telgarsky, 2020](#); [Nitanda et al., 2019](#)) which assume random Gaussian initialization $\tilde{\theta}^{(0)}$, and similar to ([Taheri & Thrampoulidis, 2023](#)) the NTK separability assumption (P3) can potentially accommodate deterministic $\theta^{(0)}$. Condition (P1) appears because we allow training both layers of the model. Specifically the $L_{2,\infty}$ norm originates from the Hessian bounds in Corollary [1](#).*

5 Application to tokenized-mixture model

We now demonstrate through an example how our results apply to specific data models.

Data model: An example. Consider $M+2$ distinct patterns $\{\boldsymbol{\mu}_+, \boldsymbol{\mu}_-, \boldsymbol{\nu}_1, \boldsymbol{\nu}_2, \dots, \boldsymbol{\nu}_M\}$, where discriminative patterns $\boldsymbol{\mu}_\pm$ correspond to labels $y = \pm 1$. The tokens are split into (i) a label-relevant set (\mathcal{R}) and (ii) a label-irrelevant set ($\mathcal{R}^c := [T] \setminus \mathcal{R}$). Conditioned on the label and \mathcal{R} , the tokens $\mathbf{x}_t, t \in [T]$ are IID as follows

$$\mathbf{x}_t|y = \begin{cases} \boldsymbol{\mu}_y & , t \in \mathcal{R} \\ \boldsymbol{\nu}_{j_t} + \mathbf{z}_t & , j_t \sim \text{Unif}(1, \dots, M) \text{ and } t \in \mathcal{R}^c, \end{cases} \quad (\text{DM1})$$

where \mathbf{z}_t are noise vectors. Let \mathcal{D} denote the joint distribution induced by the described (\mathbf{X}, y) pairs.

Assumption 3. *The labels are equi-probable and we further assume the following:*

- **Orthogonal, equal-energy means:** All patterns are orthogonal to each other, i.e. $\boldsymbol{\mu}_+ \perp \boldsymbol{\mu}_- \perp \boldsymbol{\nu}_\ell \perp \boldsymbol{\nu}_{\ell'}, \forall \ell, \ell' \in [M]$. Also, for all $y \in \{\pm 1\}, \ell \in [M]$ that $\|\boldsymbol{\mu}_y\| = \|\boldsymbol{\nu}_\ell\| = S$, where S denotes the signal strength.
- **Sparsity level:** The number of label-relevant tokens is $|\mathcal{R}| = \zeta T$; for sparsity level $\zeta \in (0, 1)$.
- **Noise distribution:** The noise tokens \mathbf{z}_t are sampled from a distribution \mathcal{D}_z , such that it holds almost surely for $\mathbf{z}_t \sim \mathcal{D}_z$ that $|\langle \mathbf{z}_t, \boldsymbol{\mu}_y \rangle| \leq Z_\mu, y \in \{\pm 1\}$ and $|\langle \mathbf{z}_t, \boldsymbol{\nu}_\ell \rangle| \leq Z_\nu/M, \forall \ell \in [M]$. Moreover, $\|\mathbf{z}_t\| \leq Z$. Overall, Assumption 1 is satisfied with $R = \sqrt{S^2 + Z^2 + 2Z_\nu/M}$.

The above assumptions can be relaxed, but without contributing new insights. We have chosen to present a model that is representative and transparent in its analysis.

5.1 Finding a good initialization

To apply the general Corollary 2 to the specific data model DM1, it suffices to find good initialization. While we cannot directly show that $\boldsymbol{\theta}^{(0)} = \mathbf{0}$ is good, we can show this is the case for first step of gradient descent $\tilde{\boldsymbol{\theta}}^{(1)}$. Thus, we consider training in two phases as follows.

First phase: One step of GD as initialization. We use n_1 training samples to update the model parameters by running one-step of gradient descent starting from zero initialization. Specifically,

$$(\mathbf{U}_h^{(1)}, \mathbf{W}_h^{(1)}) = \boldsymbol{\theta}_h^{(1)} = \boldsymbol{\theta}_h^{(0)} - \alpha_h \sqrt{H} \cdot \nabla_{\boldsymbol{\theta}_h} \widehat{L}_{n_1}(\boldsymbol{\theta}_h^{(0)}), \text{ where } \boldsymbol{\theta}_h^{(0)} = \mathbf{0} \forall h \in [H].$$

Here, α_h denotes the step-size for head $h \in [H]$ and the scaling by \sqrt{H} guarantees the update of each head is $\mathcal{O}(1)$. The lemma below shows that at the end of this phase, we have $\|\mathbf{U}_h^{(1)} - \frac{\zeta \alpha_h}{2} \mathbf{1}_T \mathbf{u}_\star^\top\|_F = \mathcal{O}(1/\sqrt{n_1})$, where \mathbf{u}_\star is the oracle classifier $\mathbf{u}_\star = \boldsymbol{\mu}_+ - \boldsymbol{\mu}_-$. On the other hand, the attention weight-matrix does not get updated; the interesting aspect of the training lies in the second phase, which involves updating \mathbf{W} .

Lemma 2 (First phase). *After the first-gradient step as described above, we have $\mathbf{U}_h^{(1)} = \alpha_h \mathbf{1}_T (\frac{\zeta}{2} \mathbf{u}_\star^\top + \mathbf{p}^\top)$ and $\mathbf{W}_h^{(1)} = \mathbf{0}$. where with probability at least $1 - \delta \in (0, 1)$ over the randomness of labels there exists positive universal constant $C > 0$ such that*

$$\|\mathbf{p}\| \leq C (2S + Z) \left(\sqrt{\frac{d}{n_1}} + \sqrt{\frac{\log(1/\delta)}{n_1}} \right) =: P. \quad (6)$$

Second phase: GD with constant step size. During the second phase, K gradient steps are performed on n new samples (distinct from those used in the first phase). Concretely, $\tilde{\boldsymbol{\theta}}^{(k+1)} = \tilde{\boldsymbol{\theta}}^{(k)} - \eta \cdot \nabla_{\tilde{\boldsymbol{\theta}}} \widehat{L}_n(\tilde{\boldsymbol{\theta}}^{(k)})$, $k = 1, \dots, K$, with $\tilde{\boldsymbol{\theta}}^{(1)} = \text{concat}(\{\boldsymbol{\theta}_h^{(1)}\}_{h \in [H]})$ the step obtained by the first-phase update and η the step-size of the second phase. In order to analyze the second phase, during which both $\tilde{\mathbf{W}}$ and $\tilde{\mathbf{U}}$ get updated, we employ the general results of Section 4. To do so, we show that $\tilde{\boldsymbol{\theta}}^{(1)}$ serves as good initialization as per Definition 1.

Proposition 2. *Consider the first-phase iterate $\{\boldsymbol{\theta}_h^{(1)}\}_{h \in [H]}$ and condition on the event $\|\mathbf{p}\| \leq P$ (depending only on the data randomness in the first phase) of Lemma 2. Suppose the step-size of the first phase is chosen IID $\alpha_h \sim \text{Unif}(\pm 1), h \in [H]$. Then, the initialization $\tilde{\boldsymbol{\theta}}^{(1)} = \text{concat}(\boldsymbol{\theta}_1^{(1)}, \dots, \boldsymbol{\theta}_H^{(1)})$ is good with respect to data sampled from DM1 and satisfying Assumption 3. Specifically, the three desired properties hold as follows.*

- Almost surely, P1 holds with $B_2 = \sqrt{T}(S + P)$.

- With probability $1 - \delta \in (0, 1)$, **P2** holds with $B_\Phi = \text{TR}(S + P)\sqrt{2\log(n/\delta)}$.
- Suppose $\sqrt{H} \gtrsim \frac{R^4 T(S+P)}{\gamma_*} \cdot \sqrt{2\log(n/\delta)}$. Then, with probability $1 - \delta \in (0, 1)$, **P3** holds with $\gamma = \gamma_*/2$ where

$$\gamma_* := \frac{T(1-\zeta)\zeta(\zeta S^4 - 7\bar{Z}S^2 - 12\bar{Z}^2 - 16\frac{\bar{Z}^3}{S^2})}{4\sqrt{2(M+1)}} - PT^{5/2}(S+Z)^3 + \frac{S\sqrt{T}\left(\zeta - 2(1-\zeta)\frac{Z_\mu}{S^2}\right)}{\sqrt{2}}, \quad (7)$$

and $\bar{Z} := Z_\mu \vee Z_\nu$. The randomness is with respect to the sampling of $\alpha_h, h \in [H]$.

The parameter γ_* in (7) represents the NTK margin of the model at initialization $\tilde{\theta}^{(1)}$. By Corollary 2, larger margin translates to better train/generalization bounds and smaller requirements on the number of heads. For a concrete example, suppose $T \vee M = \mathcal{O}(1)$ and $Z \vee \bar{Z} = \mathcal{O}(S)$. Then, provided first-phase sample size $n_1 \gtrsim S^2 d$ so that $P = \mathcal{O}(1)$, it holds $\gamma_* = \gamma_{\text{lin}} + \Omega(\zeta^2(1-\zeta)S^4)$, where $\gamma_{\text{lin}} = \Omega(\zeta S)$ is the margin of a linear model for the same dataset (see App. F). Overall, applying Cor. 2 for $K = n$ and a polylogarithmic $\text{polylog}(n)$ number of heads leads to $\tilde{\mathcal{O}}\left(\frac{1}{\eta\gamma_*^2 n}\right)$ train loss and expected generalization gap.

5.2 Proof sketch of P3: NTK separability

It is instructive to see how **P3** follows as it sheds light on the choice of an appropriate target parameter $\tilde{\theta}$ as per Thms. 1 and 2. We choose

$$\mathbf{W}_* = \mu_+ \mu_+^\top + \mu_- \mu_-^\top + \sum_{\ell \in [M]} \nu_\ell (\mu_+ + \mu_-)^\top \quad \text{and} \quad \mathbf{U}_* = \mathbf{1}_T \mathbf{u}_*^\top = \mathbf{1}_T (\mu_+ - \mu_-)^\top,$$

and normalize parameters such that $\theta_* := (\bar{\mathbf{U}}_* = \frac{1}{\|\mathbf{U}_*\|_F} \mathbf{U}_*, \text{sign}(\alpha) \bar{\mathbf{W}}_* = \text{sign}(\alpha) \frac{1}{\|\mathbf{W}_*\|_F} \mathbf{W}_*)$. It is easy to see that \mathbf{U}_* is the optimal classifier for the label-relevant tokens. To gain intuition on the choice of \mathbf{W}_* , note that $\mathbf{W}_* = \mathbf{W}_{K,*} \mathbf{W}_{Q,*}^\top$, with key-query matrices chosen as $\mathbf{W}_{K,*} = [\mu_+ \quad \mu_- \quad \nu_1 \quad \dots \quad \nu_M] \in \mathbb{R}^{d \times (M+2)}$ and $\mathbf{W}_{Q,*} = [\mu_+ \quad \mu_- \quad \mu_+ + \mu_- \quad \dots \quad \mu_+ + \mu_-] \in \mathbb{R}^{d \times (M+2)}$. With these choices, the relevance scores (aka softmax logits) of relevant tokens turn out to be strictly larger compared to the irrelevant tokens. Concretely, we show in App. D.2.3 that the t -th row $\mathbf{r}_t(\mathbf{X}; \mathbf{W}_*) = \mathbf{X} \mathbf{W}_*^\top \mathbf{x}_t$ of the softmax-logit matrix satisfies the following:

$$\forall t : [\mathbf{r}_t]_{t'} = \begin{cases} \mathcal{O}(S^4) & , t' \in \mathcal{R}, \\ \mathcal{O}(S^2) & , t' \in \mathcal{R}^c. \end{cases} \quad (8)$$

Thus, under this parameter choice, softmax can attend to label-relevant tokens and suppresses noisy irrelevant tokens. In turn, this increases the signal-to-noise ratio for classification using \mathbf{U}_* .

We now show how to compute $\mathbb{E}_{\theta^{(1)}} y \langle \nabla_{\theta} \Phi(\mathbf{X}; \theta^{(1)}), \theta_* \rangle$ for a single head. Recall θ_* consists of $\bar{\mathbf{U}}_*, \bar{\mathbf{W}}_*$. First, since $\mathbf{W}^{(1)} = \mathbf{0}$, using Assumption 3, a simple calculation shows $y \langle \nabla_{\mathbf{U}} \Phi(\mathbf{X}; \theta^{(1)}), \bar{\mathbf{U}}_* \rangle \geq \frac{S\sqrt{T}}{\sqrt{2}} \left(\zeta - 2(1-\zeta) \frac{Z_\mu}{S^2} \right)$. Second, to compute $\mathbb{E}_{\alpha \sim \text{Unif}(\pm 1)} y \langle \nabla_{\mathbf{W}} \Phi(\mathbf{X}; \theta^{(1)}), \text{sign}(\alpha) \bar{\mathbf{W}}_* \rangle$ it follows from Lemma 1 that

$$\nabla_{\mathbf{W}} \Phi(\mathbf{X}; \theta^{(1)}) = \frac{\alpha \zeta}{2} \sum_{t \in [T]} \mathbf{x}_t \mathbf{u}_*^\top \mathbf{X}^\top \varphi'(\mathbf{0}) \mathbf{X} + \alpha \sum_{t \in [T]} \mathbf{x}_t \mathbf{p}^\top \mathbf{X}^\top \varphi'(\mathbf{0}) \mathbf{X}.$$

Note the first term is dominant here since the second term can be controlled by making $\|\mathbf{p}\|_2$ small as per Lemma 2. Thus, ignoring here the second term (see Appendix D.2.3 for full calculation) $y \langle \nabla_{\mathbf{W}} \Phi(\mathbf{X}; \theta^{(1)}), \mathbf{W}_* \rangle$ is governed by the following term: $\frac{\alpha \zeta}{2} \sum_{t \in [T]} y \mathbf{u}_*^\top \mathbf{X}^\top \varphi'(\mathbf{0}) \mathbf{X} \mathbf{W}_*^\top \mathbf{x}_t = \frac{\alpha \zeta}{2} \sum_{t \in [T]} y \mathbf{u}_*^\top \mathbf{X}^\top \varphi'(\mathbf{0}) \mathbf{r}_t$. Note that $\varphi'(\mathbf{0}) = \mathbf{I} - \frac{1}{T} \mathbf{1}_T \mathbf{1}_T^\top$. To simplify the exposition here, let us focus on the identity component and leave treatment of the rank-one term to the detailed proof. The corresponding term then becomes

$$\frac{\alpha \zeta}{2} \sum_{t \in [T]} \sum_{t' \in [T]} \underbrace{(y \mathbf{u}_*^\top \mathbf{x}_{t'})}_{\text{class. logits}} \cdot \underbrace{([\mathbf{r}_t]_{t'})}_{\text{softmax logits}},$$

which involves for each output token t , the sum of products over all tokens $t' \in [T]$ of softmax logits (i.e. relevant scores $[\mathbf{r}_t]_{t'}$) and corresponding classification logits (i.e. $\mathbf{y} \mathbf{u}_*^\top \mathbf{x}_{t'}$). Note that by choice of \mathbf{u}_* and \mathbf{W}_* , both the classification and softmax logits are large from label-relevant tokens, while being small for noise tokens. Intuitively, this allows for a positive margin γ_* as stated in Proposition 2. We defer the detailed calculations to Appendix D.2.3.

In the appendix, we also detail how to yield the computation for the MHA, which builds on the calculations for the single-head attention model above. In short, we simply choose multi-head parameter $\tilde{\boldsymbol{\theta}}_*$ as $\tilde{\boldsymbol{\theta}}_* = \frac{1}{\sqrt{H}} \text{concat}(\boldsymbol{\theta}_*(\boldsymbol{\theta}_1^{(1)}), \dots, \boldsymbol{\theta}_*(\boldsymbol{\theta}_H^{(1)}))$. This guarantees that $\|\tilde{\boldsymbol{\theta}}_*\| = \sqrt{2}$ and maintains the multi-head NTK margin be at least γ_* in expectation. To complete the proof, it remains to get a high-probability version of this bound. To do this, notice that $\boldsymbol{\theta}_h^{(1)}$ are IID, hence we can apply Hoeffding’s inequality, which finally gives the desired bound on the NTK margin provided sufficient number of heads H , which controls the degree of concentration when applying Hoeffding’s inequality. See Lemmas 14 and 15 for details.

5.3 Is the NTK margin optimal?

Below, we discuss the optimality of the NTK margin γ_* . First, define set of parameters $\boldsymbol{\theta}_{\text{opt}} := (\mathbf{U}_{\text{opt}}, \mathbf{W}_{\text{opt}})$:

$$\mathbf{U}_{\text{opt}} := \frac{1}{S\sqrt{2T}} \mathbf{U}_* \quad \text{and} \quad \mathbf{W}_{\text{opt}} := \frac{1}{S^2\sqrt{2(M+1)}} \mathbf{W}_*, \quad (9)$$

normalized so that $\|\boldsymbol{\theta}_{\text{opt}}\|_F = \sqrt{2}$. Recall here the definitions of $\mathbf{U}_*, \mathbf{W}_*$ in the section above. As we already explained above, this choice of parameters guarantees that relevant tokens are assigned larger relevance and classification scores compared to irrelevant ones. Specifically about \mathbf{W}_* , we saw in Eq. (8) that it ensures a gap of $\mathcal{O}(S^2)$ between relevance scores of label-relevant and label-irrelevant tokens. Thanks to this gap, it is possible for softmax to fully attend to the label-relevant tokens by saturating the softmax. To do this, it suffices to scale-up \mathbf{W}_* by an amount $\propto 1/S^2$. This is formalized in the proposition below.

Proposition 3 (Attention expressivity for tokenized mixture model). *Consider single-head attention model. Suppose the noise level is such that $Z_\mu = Z_\nu \leq S^2/8$. For any $\epsilon > 0$, consider Γ_ϵ satisfying $\Gamma_\epsilon \geq \frac{8\sqrt{2(M+1)}}{3S^2} \log\left(\frac{\zeta^{-1}-1}{\epsilon}\right)$. Then, the attention scores corresponding to weights $\Gamma_\epsilon \cdot \mathbf{W}_{\text{opt}}$ satisfy*

$$\forall t \in [T] : 0 \leq 1 - \sum_{t' \in \mathcal{R}} \varphi_{t'}(\mathbf{x}_t^\top \Gamma_\epsilon \mathbf{W}_{\text{opt}} \mathbf{X}^T) = \sum_{t' \in \mathcal{R}^c} \varphi_{t'}(\mathbf{x}_t^\top \Gamma_\epsilon \mathbf{W}_{\text{opt}} \mathbf{X}^T) \leq \epsilon. \quad (10)$$

Thus, almost surely over data (\mathbf{X}, \mathbf{y}) generated from data model DM1 the margin of single-head attention with parameters $(\mathbf{U}_{\text{opt}}, \Gamma_\epsilon \cdot \mathbf{W}_{\text{opt}})$ satisfies

$$\mathbf{y} \Phi(\mathbf{X}; \mathbf{U}_{\text{opt}}, \Gamma_\epsilon \cdot \mathbf{W}_{\text{opt}}) \geq \gamma_{\text{attn}} := \gamma_{\text{attn}}(\epsilon) := \frac{\sqrt{T}}{\sqrt{2}S} (S^2(1-\epsilon) - 2\epsilon Z_\mu). \quad (11)$$

From Eq. (10), note that as $\epsilon \rightarrow 0$ and $\Gamma_\epsilon \rightarrow \infty$, the softmax map saturates, i.e. it approaches a hard-max map that attends only to the label-relevant tokens (\mathcal{R}) and suppress the rest (\mathcal{R}^c). As a consequence of this, Eq. (11) shows that the achieved margin approaches $\gamma_{\text{attn}} := S\sqrt{T}/\sqrt{2}$. Note this is independent of the sparsity level ζ . In particular, $\gamma_{\text{attn}} \geq \gamma_* \geq \gamma_{\text{lin}}$ and the gap increases with decreasing sparsity. See appendix for experiments and discussion regarding the margin achieved by GD for data model DM1.

Following Proposition 3, a natural question arises: Is it possible to choose “good” parameters $\tilde{\boldsymbol{\theta}} = (\tilde{\mathbf{U}}, \tilde{\mathbf{W}})$ based on the set of optimal parameters $\boldsymbol{\theta}_{\text{opt}}$? This would then yield train-loss and expected generalization-gap bounds $\tilde{\mathcal{O}}(1/(\eta\gamma_{\text{attn}}^2 n))$ after $\Theta(n)$ steps of GD starting at $\tilde{\boldsymbol{\theta}}^{(0)} = \mathbf{0}$. To investigate this question, define the following parameters for each head, aligning with the aforementioned “good” directions of Proposition 3:

$$\mathbf{U}_h := \frac{\log(n)}{\gamma_{\text{attn}}} \frac{1}{H^{1/2}} \mathbf{U}_{\text{opt}}, \quad \mathbf{W}_h := \frac{C}{H^p} \mathbf{W}_{\text{opt}},$$

for some $C > 0$, $p > 0$, and $\forall h \in [H]$. To yield the margin γ_{attn} of (10), we need that each \mathbf{W}_h has norm at least $\Gamma_\epsilon \propto 1/S^2$. Thus, we need $\|\mathbf{W}_h\| \gtrsim \frac{1}{S^2} \implies S^2 \gtrsim \frac{1}{C} \cdot H^p$. Now, in order to apply Thms. 1 and 2, the

requirement on the number of heads H in terms of distance of $\tilde{\theta}$ to $\tilde{\theta}^{(0)} = \mathbf{0}$ yields the following condition:

$$H^{1/2} \gtrsim S^5 \|\tilde{\theta}\|^3. \quad (12)$$

Note that $\|\tilde{U}\| = \frac{\log(n)}{\gamma_{\text{attn}}} \approx \frac{\log(n)}{S}$, $\|\tilde{W}\| = C \cdot H^{1/2-p}$. Hence, in computing $\|\tilde{\theta}\|$, we distinguish two cases.

First, assume that $\|\tilde{W}\| \geq \|\tilde{U}\|$ which implies that $S \gtrsim \frac{\log(n)}{C} \cdot H^{p-1/2}$ and $\|\tilde{\theta}\| \gtrsim \|\tilde{W}\| \vee \|\tilde{U}\| = C \cdot H^{1/2-p}$. Since

$$S \gtrsim \frac{1}{C^{1/2}} \cdot H^{p/2} \vee \frac{\log(n)}{C} \cdot H^{p-1/2},$$

by using Eq. (12), we get the following conditions on H :

$$H^{1/2} \gtrsim S^5 \cdot C^3 \cdot H^{3/2-3p} \gtrsim C^{1/2} \cdot H^{3/2-p/2} \vee \frac{\log^5(n)}{C^2} \cdot H^{2p-1} \implies H^{p-2} \gtrsim C \text{ and } H^{p-1/4} \lesssim \frac{C}{\log^{5/2}(n)}.$$

Combining these two gives $C \lesssim \frac{C}{\log^{5/2}(n)} \implies \log(n) \lesssim 1$, a contradiction since $n > 1$. Thus, there are no possible choices for p and C that satisfy both conditions. The case $\|\tilde{W}\| \leq \|\tilde{U}\|$ can be treated similarly leading to the same conclusion; thus, is omitted for brevity.

Intuitively, this contradiction arises because of the large $\|\mathbf{W}_h\|$ requirement to achieve margin γ_{attn} . Finally, one can ask if it is possible to resolve the contradiction by changing the scaling of normalization with respect to H in the MHA model Eq. (2), from $1/H^{1/2}$ to $1/H^c$ for $c > 0$. It can be shown via the same argument that no such value of c exists for which $\tilde{\theta}$ constructed above satisfies the overparameterization requirement $H^c \gtrsim S^5 \|\tilde{\theta}\|^3$. We thus conclude that the construction of weights in Proposition 3 does not yield a target parameter that simultaneously achieves low empirical loss and allows choosing H large enough as per (3). This triggers interesting questions for future research: Does GD converge to weights attaining margin γ_{attn} as in Proposition 3? If so, under what conditions on initialization? See also the remarks in Section 7.

6 Proof Sketch of Section 4

Throughout this section we drop the \sim in $\tilde{\theta}$ and $\tilde{\Phi}(\mathbf{X}_i; \tilde{\theta})$ as everything refers to the full model. Moreover, $\tilde{\theta}^{(K)}$ and $\tilde{\theta}^{(0)}$ are denoted by θ_K and θ_0 .

6.1 Training analysis

The proof begins by showing step-wise descent for any iteration $k \geq 0$ of GD (see Lemma 7), where step-size at each iteration $\eta_k \leq \frac{1}{\rho_k}$ depends on the objective's local smoothness parameters $\rho_k = \beta_2(\theta_k) \vee \beta_2(\theta_{k+1})$:

$$\widehat{L}(\theta_{k+1}) \leq \widehat{L}(\theta_k) - \frac{\eta_k}{2} \|\nabla \widehat{L}(\theta_k)\|^2. \quad (13)$$

Now, using Taylor's theorem we can link $\widehat{L}(\theta_k)$ to $\widehat{L}(\theta)$ for any θ as follows:

$$\widehat{L}(\theta) \leq \widehat{L}(\theta_k) + \langle \nabla \widehat{L}(\theta_k), \theta - \theta_k \rangle + \frac{1}{2} \min_{\theta_{k_\alpha}} \lambda_{\min}(\nabla^2 \widehat{L}(\theta_{k_\alpha})) \|\theta - \theta_k\|^2, \quad (14)$$

where $\theta_{k_\alpha} := \alpha \theta_k + (1 - \alpha) \theta$, $\alpha \in [0, 1]$. We can plug this into (13) to relate the loss at iterates θ_k and θ_{k+1} . To continue, we need to lower bound $\min_{\theta_{k_\alpha}} \lambda_{\min}(\nabla^2 \widehat{L}(\theta_{k_\alpha}))$. For this, we use the following property of the loss objective from Corollary 1: $\forall \theta : \lambda_{\min}(\nabla^2 \widehat{L}(\theta)) \geq -\kappa(\theta) \cdot \widehat{L}(\theta)$, where $\kappa(\theta) := \frac{\beta_3(\theta)}{\sqrt{H}}$. Note from the definition of $\beta_3(\cdot)$ that $\forall \theta_1, \theta_2 : \max_{\theta \in [\theta_1, \theta_2]} \beta_3(\theta) = \beta_3(\theta_1) \vee \beta_3(\theta_2)$. Thus, the above property of the loss implies the following *local self-bounded weak convexity* property on the line $[\theta_1, \theta_2]$ for arbitrary points θ_1, θ_2 :

$$\forall \theta_1, \theta_2 : \min_{\theta \in [\theta_1, \theta_2]} \lambda_{\min}(\nabla^2 \widehat{L}(\theta)) \geq -\frac{\beta_3(\theta_1) \vee \beta_3(\theta_2)}{\sqrt{H}} \cdot \max_{\theta \in [\theta_1, \theta_2]} \widehat{L}(\theta). \quad (15)$$

Therefore, using Eq. (15) in Eq. (14), we can get:

$$\widehat{L}(\boldsymbol{\theta}) \geq \widehat{L}(\boldsymbol{\theta}_k) + \langle \nabla \widehat{L}(\boldsymbol{\theta}_k), \boldsymbol{\theta} - \boldsymbol{\theta}_k \rangle - \frac{1}{2} \frac{\beta_3(\boldsymbol{\theta}_1) \vee \beta_3(\boldsymbol{\theta}_2)}{\sqrt{H}} \cdot \max_{\alpha \in [0,1]} \widehat{L}(\boldsymbol{\theta}_{k_\alpha}) \|\boldsymbol{\theta} - \boldsymbol{\theta}_k\|^2. \quad (16)$$

To apply the Descent Lemma in (13), we need to fix a step-size such that satisfies the condition of the Lemma at each iteration $\eta \leq \eta_k$ for all $k < K$. Then, combining with Eq. (16) and applying standard telescope summation, we arrive at the following:

$$\frac{1}{K} \sum_{k=1}^K \widehat{L}(\boldsymbol{\theta}_k) \leq \widehat{L}(\boldsymbol{\theta}) + \frac{\|\boldsymbol{\theta} - \boldsymbol{\theta}_0\|^2}{2\eta K} + \frac{1}{2K} \sum_{k=0}^{K-1} \frac{\beta_3(\boldsymbol{\theta}) \vee \beta_3(\boldsymbol{\theta}_k)}{\sqrt{H}} \cdot \max_{\alpha \in [0,1]} \widehat{L}(\boldsymbol{\theta}_{k_\alpha}) \|\boldsymbol{\theta} - \boldsymbol{\theta}_k\|^2. \quad (17)$$

Next, we use the following generalized local quasi-convexity (GLQC) of the loss function.

Proposition 4 (GLQC property: Slight variation of Prop. 8 of Taheri & Thrampoulidis (2023)). *Let $\boldsymbol{\theta}_1$ and $\boldsymbol{\theta}_2$ be two points that are sufficiently close to each other, such that*

$$2(\beta_3(\boldsymbol{\theta}_1) \vee \beta_3(\boldsymbol{\theta}_2)) \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\|^2 \leq \sqrt{H}. \quad (18)$$

Then, $\max_{\boldsymbol{\theta} \in [\boldsymbol{\theta}_1, \boldsymbol{\theta}_2]} \widehat{L}(\boldsymbol{\theta}) \leq \frac{4}{3} (\widehat{L}(\boldsymbol{\theta}_1) \vee \widehat{L}(\boldsymbol{\theta}_2))$.

Using Proposition 4 in Eq. (17) and assuming sufficiently large heads H such that $\sqrt{H} \geq 2(\beta_3(\boldsymbol{\theta}) \vee \beta_3(\boldsymbol{\theta}_k)) \|\boldsymbol{\theta} - \boldsymbol{\theta}_k\|^2$, we can get the advertised regret bound in (4).

In order to remove the dependence of H on iteration k , by an induction argument we can show bounded iterates-norm i.e. $\|\boldsymbol{\theta}_k - \boldsymbol{\theta}\| \leq 3\|\boldsymbol{\theta} - \boldsymbol{\theta}_0\|$ (see Lemma 10). Using this and the definition of $\beta_3(\cdot)$ we can control $\beta_3(\boldsymbol{\theta}) \vee \beta_3(\boldsymbol{\theta}_k)$ as $(\beta_3(\boldsymbol{\theta}) \vee \beta_3(\boldsymbol{\theta}_k)) \lesssim \|\boldsymbol{\theta} - \boldsymbol{\theta}_0\| + \|\boldsymbol{\theta}\|_{2,\infty}$ to get the desired requirement of heads $\sqrt{H} \gtrsim \|\boldsymbol{\theta}\|_{2,\infty} \|\boldsymbol{\theta} - \boldsymbol{\theta}_0\|^3$ stated in Eq. (3).

The remaining piece to guarantee descent at each step is establishing a $\rho(\boldsymbol{\theta})$ such that $\rho_k \leq \rho(\boldsymbol{\theta})$ for all $k < K$. To do this, we recall that $\rho_k = \beta_2(\boldsymbol{\theta}_k) \vee \beta_2(\boldsymbol{\theta}_{k+1})$. By definition of $\beta_2(\cdot)$ in Corollary 1, we can control $\beta_2(\boldsymbol{\theta}_k) \vee \beta_2(\boldsymbol{\theta}_{k+1})$ with controlling $\|\boldsymbol{\theta}_k\|_{2,\infty} \vee \|\boldsymbol{\theta}_{k+1}\|_{2,\infty}$ as $(\|\boldsymbol{\theta}_k\|_{2,\infty} \vee \|\boldsymbol{\theta}_{k+1}\|_{2,\infty}) \lesssim \|\boldsymbol{\theta} - \boldsymbol{\theta}_k\| + \|\boldsymbol{\theta}\|_{2,\infty} + 1$. Using iterates-norm bound and setting $\rho(\boldsymbol{\theta}) = \left(\frac{2\sqrt{T}dR^3}{\sqrt{H}} + \frac{TR^2}{4}\right)\alpha(\boldsymbol{\theta})^2$ with $\alpha(\boldsymbol{\theta}) := 3\sqrt{d}R^2[3\sqrt{T}R^3(3\|\boldsymbol{\theta} - \boldsymbol{\theta}_0\| + \|\boldsymbol{\theta}\|_{2,\infty}) + 2\sqrt{T}R]$, satisfies the desired condition for the Descent Lemma completing the proof.

6.2 Generalization analysis

In order to bound the expected generalization gap, we leverage the algorithmic stability framework. To begin, consider the leave-one-out (loo) training loss $\widehat{L}^{-i}(\boldsymbol{\theta}) := \frac{1}{n} \sum_{j \neq i} \ell_j(\boldsymbol{\theta})$ for $i \in [n]$, where $\ell_j(\boldsymbol{\theta}) := \ell(y_j \Phi(\mathbf{X}_j; \boldsymbol{\theta}))$ denotes the j -th sample loss. With these, define the loo model updates of GD on the loo loss for $\eta > 0$:

$$\boldsymbol{\theta}_{k+1}^{-i} := \boldsymbol{\theta}_k^{-i} - \eta \nabla \widehat{L}^{-i}(\boldsymbol{\theta}_k^{-i}), \quad k \geq 0, \quad \boldsymbol{\theta}_0^{-i} = \boldsymbol{\theta}_0.$$

The following lemma relates expected generalization loss to average model stability for any G -Lipschitz loss.

Lemma 3 (Lei & Ying (2020), Thm. 2). *For G -Lipschitz loss and for all iterates K , it holds that $\mathbb{E}[L(\boldsymbol{\theta}_K) - \widehat{L}(\boldsymbol{\theta}_K)] \leq 2G \cdot \mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n \|\boldsymbol{\theta}_K - \boldsymbol{\theta}_K^{-i}\|\right]$.*

To bound the average model-stability on the r.h.s of the lemma's inequality, we use GD expansiveness. Specifically applying (Taheri & Thrampoulidis, 2023, Lemma B.1.) to our setting, gives $\forall \boldsymbol{\theta}, \boldsymbol{\theta}'$:

$$\|(\boldsymbol{\theta} - \eta \nabla \widehat{L}(\boldsymbol{\theta})) - (\boldsymbol{\theta}' - \eta \nabla \widehat{L}(\boldsymbol{\theta}'))\| \leq \max_{\alpha \in [0,1]} \left\{ \left(1 + \frac{\eta \beta_3(\boldsymbol{\theta}_\alpha)}{\sqrt{H}} \widehat{L}(\boldsymbol{\theta}_\alpha)\right) \vee \eta \beta_2(\boldsymbol{\theta}_\alpha) \right\} \|\boldsymbol{\theta} - \boldsymbol{\theta}'\|, \quad (19)$$

where, $\boldsymbol{\theta}_\alpha = \alpha \boldsymbol{\theta} + (1 - \alpha) \boldsymbol{\theta}'$, $\alpha \in [0, 1]$. Using this and gradient self-boundedness from Corollary 1, we get:

$$\|\boldsymbol{\theta}_{k+1} - \boldsymbol{\theta}_{k+1}^{-i}\| \leq \max_{\alpha \in [0,1]} \left\{ \left(1 + \frac{\eta \beta_3(\boldsymbol{\theta}_{k_\alpha}^{-i})}{\sqrt{H}} \widehat{L}^{-i}(\boldsymbol{\theta}_{k_\alpha}^{-i})\right) \vee \eta \beta_2(\boldsymbol{\theta}_{k_\alpha}^{-i}) \right\} \cdot \|\boldsymbol{\theta}_k - \boldsymbol{\theta}_k^{-i}\| + \frac{\eta \beta_1(\boldsymbol{\theta}_k)}{n} \ell_i(\boldsymbol{\theta}_k), \quad (20)$$

where $\theta_{k\alpha}^{-i} := \alpha\theta_k + (1-\alpha)\theta_k^{-i}$ for $\alpha \in [0, 1]$. Further using the bounded iterates-norm property from the training analysis, we can control $\beta_2(\theta_{k\alpha}^{-i}) \leq \tilde{\beta}_2(\theta)$ and $\beta_3(\theta_{k\alpha}^{-i}) \leq \tilde{\beta}_3(\theta)$ making them independent of k (See Lemma 11 for the definitions of $\tilde{\beta}_2(\cdot), \tilde{\beta}_3(\cdot)$). In order to invoke the Descent Lemma, we set the step-size same as in the training analysis. Thus, (20) becomes:

$$\|\theta_{k+1} - \theta_{k+1}^{-i}\| \leq \left(1 + \frac{\eta\tilde{\beta}_3(\theta)}{\sqrt{H}}\right) \max_{\alpha \in [0,1]} \widehat{L}^{-i}(\theta_{k\alpha}^{-i}) \|\theta_k - \theta_k^{-i}\| + \frac{\eta\beta_1(\theta_k)}{n} \ell_i(\theta_k). \quad (21)$$

As in the training analysis, we can control the loo empirical loss \widehat{L}^{-i} for any point on the line $[\theta_k, \theta_k^{-i}]$ of two sufficiently close points satisfying $\sqrt{H} \geq 2(\beta_3(\theta_k) \vee \beta_3(\theta_k^{-i})) \|\theta_k - \theta_k^{-i}\|^2$. Using Prop. 4, Eq. (21) becomes

$$\|\theta_{k+1} - \theta_{k+1}^{-i}\| \leq (1 + \alpha_{k,i}) \|\theta_k - \theta_k^{-i}\| + \frac{\eta\tilde{\beta}_1(\theta)}{n} \ell_i(\theta_k), \quad (22)$$

where $\alpha_{k,i} := \frac{4\eta\tilde{\beta}_3(\theta)}{3\sqrt{H}} (\widehat{L}^{-i}(\theta_k) + \widehat{L}^{-i}(\theta_k^{-i}))$ and $\beta_1(\theta_k) \leq \tilde{\beta}_1(\theta)$ similar to $\beta_2(\cdot), \beta_3(\cdot)$ using bounded iterates-norm. Unrolling the iterates in (22), summing over $i \in [n]$ and using training regret bounds, we have the following average model stability bound for any iterate K : $\frac{1}{n} \sum_{i=1}^n \|\theta_K - \theta_K^{-i}\| \leq \frac{2\eta\tilde{\beta}_1(\theta)}{n} (2K\widehat{L}(\theta) + \frac{9\|\theta - \theta_0\|^2}{4\eta})$. Combining this with an application of Lemma 3 for our objective, which is $G \leq \tilde{\beta}_1(\theta)$ -Lipschitz from Corollary 1, and using $\eta \leq \frac{1}{\rho(\theta)} \leq \frac{1}{(\tilde{\beta}_1(\theta))^2}$, we get the desired generalization gap stated in Thm. 2.

7 Concluding remarks

We studied convergence and generalization of GD for training a multi-head attention layer in a classification task. Our training and generalization bounds hold under an appropriate realizability condition asking for the existence of an a target model $\tilde{\theta}$ achieving good train loss while being sufficiently close to initialization. In particular, from the condition on the number of heads H in (3), we need $\tilde{\theta}$ is at most $\tilde{\mathcal{O}}(d^{-1/3}T^{-1/6}R^{-5/3}H^{1/6})$ far from initialization (provided $\|\tilde{\theta}\|_{2,\infty} = \mathcal{O}(1)$). In Sec. 4.3 we showed that such a model exists if the initialization is chosen appropriately. Specifically it suffices that $\|\tilde{\theta}^{(0)}\|_{2,\infty} = \mathcal{O}(1)$, the model output at initialization is $\tilde{\mathcal{O}}(1)$ -bounded and that the data are linearly separable with margin γ with respect to the NTK features of the model at initialization. Then, $\mathcal{O}(d^2TR^{10}\text{polylog}(n)/\gamma^6)$ number of heads guarantee that $\Theta(n)$ GD steps result in train and test loss bounds $\tilde{\mathcal{O}}(1/(\eta\gamma^2n))$. In Sec. 5 we applied our results to a tokenized-mixture model. We showed that after one randomized gradient step from $\mathbf{0}$, the model satisfies the above conditions for good initialization. For this initialization, we computed the NTK margin γ_* which in turn governs the guaranteed rate of convergence and generalization based on our general bounds. This opens several interesting questions for future work.

First, does random initialization of attention weights satisfy NTK separability, and if so, what is the corresponding margin? Second, are there other initialization strategies that guarantee the realizability conditions are satisfied? Here, note that our conditions for good initialization are only shown to be sufficient for realizability leaving room for improvements. Third, how suboptimal is the best NTK margin (among other potential natural initializations) compared to the model's global margin $\arg \max_{\|\tilde{\theta}\|=\sqrt{2}} \min_{i \in [n]} y_i \tilde{\Phi}(\mathbf{X}_i; \tilde{\theta})$? In Proposition 3 we showed for the data model DMI that there exists single-head attention model $\theta_{\text{opt}} = (\mathbf{U}_{\text{opt}}, \mathbf{W}_{\text{opt}})$ with $\|\theta_{\text{opt}}\| = \sqrt{2}$ such that $y\Phi(\mathbf{X}; \mathbf{U}_{\text{opt}}, \Gamma_\epsilon \cdot \mathbf{W}_{\text{opt}}) = \frac{S\sqrt{T}}{\sqrt{2}} ((1-\epsilon) - 2\epsilon Z_\mu/S)$ for all $\Gamma_\epsilon \gtrsim \frac{\log((\zeta^{-1}-1)/\epsilon)}{S}$ and any $\epsilon \in (0, 1)$ (see App. E). In particular, as $\epsilon \rightarrow 0$ and $\Gamma_\epsilon \rightarrow \infty$ (for which the softmax map gets saturated and attends to tokens with highest relevance score) the achieved margin approaches $\gamma_{\text{attn}} := S\sqrt{T}/\sqrt{2}$, which is independent of the sparsity level ζ . In particular, $\gamma_{\text{attn}} \geq \gamma_* \geq \gamma_{\text{lin}}$ and the gap increases with decreasing sparsity. Is it possible to establish finite-time convergence bounds to models with margin $\approx \gamma_{\text{attn}}$ under appropriate initialization? How is the answer affected by the fact that the optimal attention weights in this case are diverging in norm ($\Gamma_\epsilon \rightarrow \infty$)? Using our approach, we argued in Sec. 5.3 that the key challenge is the saturation of norm of $\mathbf{W}_{\text{opt}}(\Gamma_\epsilon)$, which does not allow the appropriate realizability condition to hold (at least for $\mathbf{0}$ initialization). Finally, it is interesting to consider other data models for which multiple heads are necessary to interpolate the data.

References

- Ekin Akyürek, Dale Schuurmans, Jacob Andreas, Tengyu Ma, and Denny Zhou. What learning algorithm is in-context learning? investigations with linear models. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=0g0X4H8yN4I>.
- Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song. A convergence theory for deep learning via overparameterization. In *International Conference on Machine Learning*, pp. 242–252. PMLR, 2019.
- Sanjeev Arora, Simon Du, Wei Hu, Zhiyuan Li, and Ruosong Wang. Fine-grained analysis of optimization and generalization for overparameterized two-layer neural networks. In *International Conference on Machine Learning*, pp. 322–332. PMLR, 2019.
- Pierre Baldi and Roman Vershynin. The quarks of attention. *arXiv preprint arXiv:2202.08371*, 2022.
- Arindam Banerjee, Pedro Cisneros-Velarde, Libin Zhu, and Mikhail Belkin. Restricted strong convexity of deep learning models with smooth activations. *arXiv preprint arXiv:2209.15106*, 2022.
- Alberto Bietti, Vivien Cabannes, Diane Bouchacourt, Herve Jegou, and Leon Bottou. Birth of a transformer: A memory viewpoint. *arXiv preprint arXiv:2306.00802*, 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 H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 1877–1901. Curran Associates, Inc., 2020. URL https://proceedings.neurips.cc/paper_files/paper/2020/file/1457c0d6bfc4967418bfb8ac142f64a-Paper.pdf.
- Yuan Cao and Quanquan Gu. Generalization bounds of stochastic gradient descent for wide and deep neural networks. *Advances in neural information processing systems*, 32, 2019.
- Zixiang Chen, Yuan Cao, Difan Zou, and Quanquan Gu. How much over-parameterization is sufficient to learn deep relu networks? In *International Conference on Learning Representations*, 2020.
- Jianpeng Cheng, Li Dong, and Mirella Lapata. Long short-term memory-networks for machine reading. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pp. 551–561. Association for Computational Linguistics, November 2016. doi: 10.18653/v1/D16-1053. URL <https://aclanthology.org/D16-1053>.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. pp. 4171–4186, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1423. URL <https://aclanthology.org/N19-1423>.
- Yihe Dong, Jean-Baptiste Cordonnier, and Andreas Loukas. Attention is not all you need: Pure attention loses rank doubly exponentially with depth. In *International Conference on Machine Learning*, pp. 2793–2803. PMLR, 2021.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale, 2021.
- Simon Du, Jason Lee, Haochuan Li, Liwei Wang, and Xiyu Zhai. Gradient descent finds global minima of deep neural networks. In *International conference on machine learning*, pp. 1675–1685. PMLR, 2019.
- Benjamin L Edelman, Surbhi Goel, Sham Kakade, and Cyril Zhang. Inductive biases and variable creation in self-attention mechanisms. *arXiv preprint arXiv:2110.10090*, 2021.

- Tolga Ergen, Behnam Neyshabur, and Harsh Mehta. Convexifying transformers: Improving optimization and understanding of transformer networks. *arXiv preprint arXiv:2211.11052*, 2022.
- Jiri Hron, Yasaman Bahri, Jascha Sohl-Dickstein, and Roman Novak. Infinite attention: Nngp and ntk for deep attention networks. In *International Conference on Machine Learning*, pp. 4376–4386. PMLR, 2020.
- Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization in neural networks. *Advances in neural information processing systems*, 31, 2018.
- Samy Jelassi, Michael Eli Sander, and Yuanzhi Li. Vision transformers provably learn spatial structure. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho (eds.), *Advances in Neural Information Processing Systems*, 2022. URL <https://openreview.net/forum?id=eMW9AkXaREI>.
- Ziwei Ji and Matus Telgarsky. Risk and parameter convergence of logistic regression. *arXiv preprint arXiv:1803.07300*, 2018.
- Ziwei Ji and Matus Telgarsky. Polylogarithmic width suffices for gradient descent to achieve arbitrarily small test error with shallow relu networks. In *International Conference on Learning Representations*, 2020.
- Ziwei Ji and Matus Telgarsky. Characterizing the implicit bias via a primal-dual analysis. In *Algorithmic Learning Theory*, pp. 772–804. PMLR, 2021.
- Yunwen Lei and Yiming Ying. Fine-grained analysis of stability and generalization for stochastic gradient descent. In *International Conference on Machine Learning*, pp. 5809–5819. PMLR, 2020.
- Yunwen Lei, Rong Jin, and Yiming Ying. Stability and generalization analysis of gradient methods for shallow neural networks. In *Advances in Neural Information Processing Systems*, 2022.
- Hongkang Li, Meng Weng, Sijia Liu, and Pin-Yu Chen. A theoretical understanding of shallow vision transformers: Learning, generalization, and sample complexity. In *International Conference on Learning Representations*, 2023a.
- Yingcong Li, M. Emrullah Ildiz, Dimitris Papailiopoulos, and Samet Oymak. Transformers as algorithms: Generalization and stability in in-context learning, 2023b.
- Valerii Likhoshesterov, Krzysztof Choromanski, and Adrian Weller. On the expressive power of self-attention matrices, 2021.
- Zhouhan Lin, Minwei Feng, Cicero Nogueira dos Santos, Mo Yu, Bing Xiang, Bowen Zhou, and Yoshua Bengio. A structured self-attentive sentence embedding. In *International Conference on Learning Representations*, 2017. URL https://openreview.net/forum?id=BJC_jUqxe.
- Sadeh Mahdavi, Renjie Liao, and Christos Thrampoulidis. Memorization capacity of multi-head attention in transformers. *arXiv preprint arXiv:2306.02010*, 2023.
- Quynh Nguyen, Marco Mondelli, and Guido F Montufar. Tight bounds on the smallest eigenvalue of the neural tangent kernel for deep relu networks. In Marina Meila and Tong Zhang (eds.), *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pp. 8119–8129. PMLR, 18–24 Jul 2021.
- Quynh N Nguyen and Marco Mondelli. Global convergence of deep networks with one wide layer followed by pyramidal topology. *Advances in Neural Information Processing Systems*, 33:11961–11972, 2020.
- Konstantinos E Nikolakakis, Farzin Haddadpour, Amin Karbasi, and Dionysios S Kalogerias. Beyond Lipschitz: Sharp generalization and excess risk bounds for full-batch gd. *arXiv preprint arXiv:2204.12446*, 2022.
- Atsushi Nitanda, Geoffrey Chinot, and Taiji Suzuki. Gradient descent can learn less over-parameterized two-layer neural networks on classification problems. *arXiv preprint arXiv:1905.09870*, 2019.
- OpenAI. Openai: Introducing chatgpt, 2022. URL <https://openai.com/blog/chatgpt>, 2022.

OpenAI. Gpt-4 technical report, 2023.

Samet Oymak and Mahdi Soltanolkotabi. Toward moderate overparameterization: Global convergence guarantees for training shallow neural networks. *IEEE Journal on Selected Areas in Information Theory*, 1(1):84–105, 2020.

Samet Oymak, Ankit Singh Rawat, Mahdi Soltanolkotabi, and Christos Thrampoulidis. On the role of attention in prompt-tuning. In *ICLR 2023 Workshop on Mathematical and Empirical Understanding of Foundation Models*, 2023.

Ankur Parikh, Oscar Täckström, Dipanjan Das, and Jakob Uszkoreit. A decomposable attention model for natural language inference. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pp. 2249–2255, Austin, Texas, November 2016. Association for Computational Linguistics. doi: 10.18653/v1/D16-1244. URL <https://aclanthology.org/D16-1244>.

Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. In Marina Meila and Tong Zhang (eds.), *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pp. 8748–8763. PMLR, 18–24 Jul 2021. URL <https://proceedings.mlr.press/v139/radford21a.html>.

Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. 21(1), 2020. ISSN 1532-4435.

Dominic Richards and Ilja Kuzborskij. Stability & generalisation of gradient descent for shallow neural networks without the neural tangent kernel. *Advances in Neural Information Processing Systems*, 34: 8609–8621, 2021.

Dominic Richards and Mike Rabbat. Learning with gradient descent and weakly convex losses. In *International Conference on Artificial Intelligence and Statistics*, pp. 1990–1998. PMLR, 2021.

Itay M Safran, Gilad Yehudai, and Ohad Shamir. The effects of mild over-parameterization on the optimization landscape of shallow relu neural networks. In Mikhail Belkin and Samory Kpotufe (eds.), *Proceedings of Thirty Fourth Conference on Learning Theory*, volume 134 of *Proceedings of Machine Learning Research*, pp. 3889–3934. PMLR, 15–19 Aug 2021.

Arda Sahiner, Tolga Ergen, Batu Ozturkler, John Pauly, Morteza Mardani, and Mert Pilanci. Unraveling attention via convex duality: Analysis and interpretations of vision transformers. *International Conference on Machine Learning*, 2022.

Clayton Sanford, Daniel Hsu, and Matus Telgarsky. Representational strengths and limitations of transformers. *arXiv preprint arXiv:2306.02896*, 2023.

Matan Schliserman and Tomer Koren. Stability vs implicit bias of gradient methods on separable data and beyond. In Po-Ling Loh and Maxim Raginsky (eds.), *Proceedings of Thirty Fifth Conference on Learning Theory*, volume 178 of *Proceedings of Machine Learning Research*, pp. 3380–3394. PMLR, 02–05 Jul 2022.

Ohad Shamir. Gradient methods never overfit on separable data. *Journal of Machine Learning Research*, 22(85):1–20, 2021.

Daniel Soudry, Elad Hoffer, Mor Shpigel Nacson, Suriya Gunasekar, and Nathan Srebro. The implicit bias of gradient descent on separable data. *The Journal of Machine Learning Research*, 19(1):2822–2878, 2018.

Hossein Taheri and Christos Thrampoulidis. Generalization and stability of interpolating neural networks with minimal width. *arXiv preprint arXiv:2302.09235*, 2023.

Davoud Ataee Tarzanagh, Yingcong Li, Christos Thrampoulidis, and Samet Oymak. Transformers as support vector machines, 2023a.

- Davoud Ataee Tarzanagh, Yingcong Li, Xuechen Zhang, and Samet Oymak. Max-margin token selection in attention mechanism, 2023b.
- Matus Telgarsky. Margins, shrinkage, and boosting. In *International Conference on Machine Learning*, pp. 307–315. PMLR, 2013.
- Matus Telgarsky. Feature selection and low test error in shallow low-rotation relu networks. In *The Eleventh International Conference on Learning Representations*, 2022.
- Yuandong Tian, Yiping Wang, Beidi Chen, and Simon Du. Scan and snap: Understanding training dynamics and token composition in 1-layer transformer, 2023.
- Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Herve Jegou. Training data-efficient image transformers & distillation through attention. In Marina Meila and Tong Zhang (eds.), *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pp. 10347–10357. PMLR, 18–24 Jul 2021. URL <https://proceedings.mlr.press/v139/touvron21a.html>.
- 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.
- Ashish Vaswani, Noam M. Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NIPS*, 2017.
- Johannes von Oswald, Eyvind Niklasson, E. Randazzo, João Sacramento, Alexander Mordvintsev, Andrey Zhmoginov, and Max Vladymyrov. Transformers learn in-context by gradient descent. *ArXiv*, abs/2212.07677, 2022.
- Weihang Xu and Simon Du. Over-parameterization exponentially slows down gradient descent for learning a single neuron. In Gergely Neu and Lorenzo Rosasco (eds.), *Proceedings of Thirty Sixth Conference on Learning Theory*, volume 195 of *Proceedings of Machine Learning Research*, pp. 1155–1198. PMLR, 2023.
- Chulhee Yun, Srinadh Bhojanapalli, Ankit Singh Rawat, Sashank J. Reddi, and Sanjiv Kumar. Are transformers universal approximators of sequence-to-sequence functions?, 2020a.
- Chulhee Yun, Yin-Wen Chang, Srinadh Bhojanapalli, Ankit Singh Rawat, Sashank J. Reddi, and Sanjiv Kumar. $o(n)$ connections are expressive enough: Universal approximability of sparse transformers, 2020b.
- Ruiqi Zhang, Spencer Frei, and Peter L. Bartlett. Trained transformers learn linear models in-context, 2023.
- Mo Zhou, Rong Ge, and Chi Jin. A local convergence theory for mildly over-parameterized two-layer neural network. In Mikhail Belkin and Samory Kpotufe (eds.), *Proceedings of Thirty Fourth Conference on Learning Theory*, volume 134 of *Proceedings of Machine Learning Research*, pp. 4577–4632. PMLR, 15–19 Aug 2021.
- Zhenyu Zhu, Fanghui Liu, Grigorios Chrysos, Francesco Locatello, and Volkan Cevher. Benign overfitting in deep neural networks under lazy training. In *International Conference on Machine Learning*, pp. 43105–43128. PMLR, 2023.