
Transformers learn through gradual rank increase

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We identify incremental learning dynamics in transformers, where the difference
2 between trained and initial weights progressively increases in rank. We rigorously
3 prove this occurs under the simplifying assumptions of diagonal weight matrices
4 and small initialization. Our experiments support the theory and also show that
5 phenomenon can occur in practice without the simplifying assumptions.

6 1 Introduction

7 The transformer architecture achieves state of the art performance in various domains, yet we still
8 lack a solid theoretical understanding of its training dynamics (Vaswani et al., 2017; Devlin et al.,
9 2019; Liu et al., 2019; Dosovitskiy et al., 2020). Nevertheless, the theoretical toolbox has matured
10 over the last years and there are promising new approaches. One important line of work examines the
11 role that initialization scale plays on the trajectory taken by gradient descent (Jacot et al., 2018; Chizat
12 et al., 2018; Geiger et al., 2019; Moroshko et al., 2020; Jacot et al., 2021; Stöger & Soltanolkotabi,
13 2021; Kim & Chung, 2022). When the weights are initialized small, it has been shown for simple
14 networks that an *incremental learning* behaviour occurs, where functions of increasing complexity
15 are learned in stages. This regime is known to be richer than the large-initialization regime¹, but the
16 incremental learning dynamics are difficult to analyze, and are so far understood only for extremely
17 simple architectures. Can we apply this analysis to transformers? Namely:

18 *Are there incremental learning dynamics when training a transformer architecture?*

19 An obstacle is that past work on incremental learning has mainly studied linear networks (Berthier,
20 2022; Arora et al., 2019; Milanese et al., 2021; Li et al., 2020; Woodworth et al., 2019; Jacot et al.,
21 2021; Gissin et al., 2019), with one paper studying nonlinear 2-layer fully-connected networks
22 (Boursier et al., 2022). In contrast, transformers have nonlinear attention heads that do not fall under
23 previous analyses: given $\mathbf{X} \in \mathbb{R}^{n \times d}$, an attention head computes

$$\text{attention}(\mathbf{X}; \mathbf{W}_K, \mathbf{W}_Q, \mathbf{W}_V, \mathbf{W}_O) = \text{smax}(\mathbf{X} \mathbf{W}_K \mathbf{W}_Q^\top \mathbf{X}^\top) \mathbf{X} \mathbf{W}_V \mathbf{W}_O^\top \quad (1)$$

24 where $\mathbf{W}_K, \mathbf{W}_Q, \mathbf{W}_V, \mathbf{W}_O \in \mathbb{R}^{d \times d'}$ are trainable matrices, and the softmax is applied row-wise. A
25 transformer is even more complex, since it is formed by stacking alternating layers of attention heads
26 and feedforward networks, along with residual connections.

27 **Main finding** Our main finding is that transformers exhibit incremental learning dynamics, where
28 *the difference between the trained and initial weights incrementally increases in rank*. Our results
29 have a theoretical component and an experimental component.

¹In the large-initialization regime, deep learning behaves as a kernel method Jacot et al. (2018); Chizat et al. (2018). Various separations with kernels are known for smaller initialization: e.g., Ghorbani et al. (2019); Abbe et al. (2022); Malach et al. (2021).

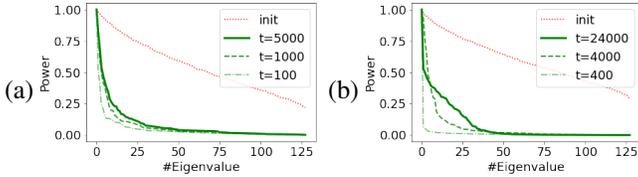


Figure 1: For an attention head in ViT trained on (a) CIFAR-10, and (b) ImageNet, we plot the normalized spectra of $\mathbf{W}_K \mathbf{W}_Q^\top$ at initialization (in red), and of the learned perturbations to $\mathbf{W}_K \mathbf{W}_Q^\top$ at different epochs (in green).

30 **Theoretical contributions** For our theory, we study a simplification of the transformer architecture, where the attention head weights are diagonal matrices: i.e., in each attention head we have
 31 $\mathbf{W}_K = \text{diag}(\mathbf{w}_K)$, where $\mathbf{w}_K \in \mathbb{R}^d$ are trainable weights, and similarly for \mathbf{W}_Q , \mathbf{W}_V and \mathbf{W}_O .
 32 We rigorously establish the training dynamics of this architecture under gradient flow when the
 33 initialization is small. We prove that dynamics occur in discrete stages: (1) during most of each stage,
 34 the loss plateaus because the weights remain close to a saddle point, and (2) at the end, the saddle
 35 point is quickly escaped and the rank of the weights increases by at most one.
 36

37 This theoretical result on transformers follows from a general theorem characterizing the learning
 38 dynamics of networks f_{NN} that depend on the product of parameters $\mathbf{u}, \mathbf{v} \in \mathbb{R}^p$ as

$$f_{\text{NN}}(\mathbf{x}; \mathbf{u}, \mathbf{v}) = h(\mathbf{x}; \mathbf{u} \odot \mathbf{v}), \quad (2)$$

39 where \mathbf{x} is the input, \odot denotes the elementwise product, and h is a smooth function.

40 **Theorem 1.1** (Informal statement of incremental learning dynamics). *Let f_{NN} be a network of*
 41 *the form (2), and suppose that the weights are initialized very small: i.e., the entries of \mathbf{u}, \mathbf{v} are*
 42 *initialized on the order $\Theta(\alpha)$ for some small $\alpha > 0$. Then the dynamics of gradient flow training*
 43 *effectively proceeds in discrete stages, each one lasting time $\Theta(\log(1/\alpha))$. In each stage, the number*
 44 *of nonnegligible entries of $\mathbf{u} \odot \mathbf{v}$ increases by at most one.*

45 A transformer with diagonal weight matrices falls under this result when we only train the attention
 46 head weights. For example, if the transformer has one attention head, then we can take $\mathbf{u} =$
 47 $[\mathbf{w}_K, \mathbf{w}_V] \in \mathbb{R}^{2d}$ and $\mathbf{v} = [\mathbf{w}_Q, \mathbf{w}_O] \in \mathbb{R}^{2d}$ to be concatenations of the diagonal entries of the
 48 weights of the head; see Example 3.2 for more details and the extension to transformers with many
 49 heads. Then, using Theorem 1.1, we see that in each stage either $\mathbf{W}_K \mathbf{W}_Q^\top = \text{diag}(\mathbf{w}_K) \text{diag}(\mathbf{w}_Q)$
 50 or $\mathbf{W}_V \mathbf{W}_O^\top = \text{diag}(\mathbf{w}_V) \text{diag}(\mathbf{w}_O)$ increases in effective rank by at most one.²

51 **Experimental contributions** In our experiments, we first validate our theoretical results, which
 52 require the simplifying assumptions of small initialization and diagonal weight matrices.

53 Then, we conduct experiments on vision transformers in settings closer to practice, without any of the
 54 assumptions required by our theoretical analysis. Perhaps surprisingly, we again observe incremental
 55 learning dynamics, even though the assumptions of the theory are not met. We observe that the
 56 difference between trained and initial weights has low rank, and also that the rank of this difference
 57 grows gradually during training; see Figure 1. The incremental nature of the dynamics is easier to see
 58 for ImageNet, since for CIFAR-10 the rank of the weight difference does not grow as much.

59 1.1 Related work

60 **Relation to LoRA** We note an intriguing connection to the LoRA algorithm, where a pretrained
 61 base model is cheaply fine-tuned by training a low-rank perturbation of the weights (Li et al., 2018;
 62 Aghajanyan et al., 2020; Hu et al., 2021). The method is surprisingly powerful, and recently LoRA
 63 has been fundamental to allowing the open-source community to inexpensively fine-tune language
 64 models (Patel & Ahmad, 2023; Taori et al., 2023). On the other hand, in our work we observe that
 65 the trained weights are a low-rank perturbation of the initial weights due to the training dynamics,
 66 without having to apply an explicit rank constraint as in LoRA. This raises an exciting open question
 67 for future work: *can we explain and improve algorithms like LoRA by better understanding and*
 68 *quantifying the incremental dynamics of large transformers?*

²We also remark that Theorem 1.1 is interesting in its own right and may have other applications beyond transformers. In fact, it qualitatively recovers the incremental dynamics result of Berthier (2022) when specialized to linear diagonal networks, i.e., when $f_{\text{NN}}(\mathbf{x}; \mathbf{u}, \mathbf{v}) = \sum_{i=1}^p u_i v_i x_i$. Furthermore, it addresses an open question of Berthier (2022) for proving incremental learning dynamics without assuming $\mathbf{u} = \mathbf{v}$ at initialization.

69 **Low-rank bias in nonlinear models** For 2-layer networks, it is known that low-rank bias in the
70 weights emerges if the target function depends on a low-dimensional subspace of the input (Abbe
71 et al., 2022, 2023; Damian et al., 2022; Bietti et al., 2022; Mousavi-Hosseini et al., 2022). The results
72 of Abbe et al. (2022, 2023) are especially relevant, since they show that the rank of the weights
73 increases in a sequential manner, determined by the “leap complexity” of the target function, which
74 is reminiscent of our empirical observations on transformers. See also Frei et al. (2022); Timor et al.
75 (2023) for more investigations of low-rank bias in 2-layer networks under different assumptions. For
76 transformers, Yu & Wu (2023) report that empirically the trained weights (using default initialization)
77 are not low-rank. This is consistent with our claim that the difference between initial and trained
78 weights is low-rank, since the initial weights might not be low-rank.

79 **Incremental learning dynamics** Several works prove incremental learning behaviour in deep
80 *linear* networks when the initialization is small. Gidel et al. (2019) has shown that gradient descent
81 dynamics on a 2-layer linear network with L_2 loss effectively solve a reduced-rank regression
82 problem with gradually increasing rank. Gissin et al. (2019) prove a dynamical depth separation
83 result, allowing for milder assumptions on initialization scale. Arora et al. (2019); Milanese et al.
84 (2021) show implicit bias towards low rank in deep matrix and tensor factorization. Li et al. (2020)
85 show deep matrix factorization dynamics with small initialization are equivalent to a greedy low-rank
86 learning (GLRL) algorithm. And Jacot et al. (2021) independently provides a similar description of
87 the dynamics, but without requiring balanced initialization. Finally, Berthier (2022); Jin et al. (2023)
88 overcome a technical hurdle from previous analyses by proving incremental learning for the entire
89 training trajectory, rather than just the first stage. In contrast to our result, these prior works apply
90 only to *linear* networks with certain convex losses, whereas our result applies to *nonlinear* networks.
91 In order to make our extension to nonlinear networks possible, we must make stronger assumptions
92 on the training trajectory, which we verify hold empirically. As far as we are aware, one other work
93 on incremental learning handles nonlinear networks: Boursier et al. (2022) proves that a 2-layer
94 network learns with a two-stage incremental dynamic; but that result needs the stylized assumption
95 that all data points are orthogonal.

96 1.2 Paper organization

97 Sections 2, 3, and 4 contain theoretical preliminaries, definitions of the models to which our theory
98 applies, and our main theoretical result on incremental dynamics. Section 5 provides experiments
99 which verify and extend the theory. Section 6 discusses limitations and future directions.

100 2 Preliminaries

101 We consider training a network $f_{\text{NN}}(\cdot; \theta)$ parametrized by a vector of weights θ , to minimize a loss

$$\mathcal{L}(\theta) = \mathbb{E}_{\mathbf{x}, \mathbf{y}}[\ell(\mathbf{y}, f_{\text{NN}}(\mathbf{x}; \theta))],$$

102 where the expectation is over samples $(\mathbf{x}, \mathbf{y}) \in \mathbb{R}^{d_y} \times \mathbb{R}^{d_x}$ from a training data distribution, and
103 $\ell : \mathbb{R}^{d_y} \times \mathbb{R}^{d_{out}} \rightarrow \mathbb{R}$. Consider a solution $\theta(t)$ to the gradient flow

$$\theta(0) = \alpha \theta_0, \quad \frac{d\theta}{dt} = -\nabla_{\theta} \mathcal{L}(\theta) \quad (3)$$

104 where $\alpha > 0$ is a parameter governing the initialization scale, that we will take very small. For our
105 theory, we henceforth require the following mild regularity assumption on the loss and data.

106 **Assumption 2.1** (Regularity of data distribution and loss). The function $\ell(\mathbf{y}, \zeta)$ is continuously
107 twice-differentiable in the arguments $[\mathbf{y}, \zeta] \in \mathbb{R}^{d_y + d_{out}}$. There exists $C > 0$ such that almost surely
108 the data is bounded by $\|\mathbf{x}\|, \|\mathbf{y}\| \leq C$.

109 The assumption on ℓ is satisfied in typical cases such as the square and the cross-entropy losses. The
110 data boundedness is often satisfied in practice (e.g., if the data is normalized).

111 3 Neural networks with diagonal weights

112 Our theory analyzes the training dynamics of networks that depend on products of diagonal weight
113 matrices. We use \odot to denote elementwise vector product.

114 **Definition 3.1.** A network f_{NN} is smooth with diagonal weights $\theta = (\mathbf{u}, \mathbf{v}) \in \mathbb{R}^{2p}$ if it is of the form

$$f_{\text{NN}}(\mathbf{x}; \theta) = h(\mathbf{x}; \mathbf{u} \odot \mathbf{v})$$

115 where $h : \mathbb{R}^{d_x} \times \mathbb{R}^p \rightarrow \mathbb{R}^{d_{out}}$ is continuously twice-differentiable in its arguments in \mathbb{R}^{d_x+p} .

116 The assumption on h precludes the use of the ReLU function since it is not continuously-differentiable.
 117 Otherwise the assumption is fairly mild since any h can be used to express an architecture of any
 118 depth as long as the nonlinearities are twice-differentiable, which includes for example GeLUs (as
 119 used in ViT). We describe how to express a transformer with diagonal weights.

120 **Example 3.2** (Transformer with diagonal weights). Consider a transformer with L layers and H
 121 attention heads on each layer. The transformer output at layer ℓ is $\mathbf{Z}_\ell \in \mathbb{R}^{n \times d}$, which is given by
 122 $\mathbf{Z}_0 = \mathbf{X}$ and inductively for $\ell > 0$ by

123 • (Attention layer) $\tilde{\mathbf{Z}}_\ell = \mathbf{Z}_{\ell-1} + \sum_{i=1}^H \text{attention}(\mathbf{Z}_{\ell-1}; \mathbf{W}_K^{\ell,i}, \mathbf{W}_Q^{\ell,i}, \mathbf{W}_V^{\ell,i}, \mathbf{W}_O^{\ell,i})$

124 • (Feedforward layer) $\mathbf{Z}_\ell = \tilde{\mathbf{Z}}_\ell + \sigma(\tilde{\mathbf{Z}}_\ell \mathbf{W}_A^\ell)(\mathbf{W}_B^\ell)^\top$,

125 where $\mathbf{W}_K^{\ell,i}, \mathbf{W}_Q^{\ell,i}, \mathbf{W}_V^{\ell,i}, \mathbf{W}_O^{\ell,i} \in \mathbb{R}^{d \times d'}$ are attention parameters, and $\mathbf{W}_A^\ell, \mathbf{W}_B^\ell \in \mathbb{R}^{d \times d'}$ are the
 126 feedforward parameters, and σ is a continuously twice-differentiable activation.

127 Suppose that the only trainable parameters are the attention parameters, and that these are diagonal
 128 matrices: i.e., $\mathbf{W}_K^{\ell,i} = \text{diag}(\mathbf{w}_K^{\ell,i})$ for some $\mathbf{w}_K^{\ell,i} \in \mathbb{R}^d$, and similarly for the other attention
 129 parameters. Because of the structure of the attention head (1), the final output \mathbf{Z}_L only depends on
 130 the attention parameters through the elementwise products $\mathbf{w}_K^{\ell,i} \odot \mathbf{w}_Q^{\ell,i}$ and $\mathbf{w}_V^{\ell,i} \odot \mathbf{w}_O^{\ell,i}$. In other
 131 words, we can write

$$\mathbf{Z}_L = h(\mathbf{X}; \mathbf{u} \odot \mathbf{v}),$$

132 for vectors $\mathbf{u} = [\mathbf{w}_K^{\ell,i}, \mathbf{w}_V^{\ell,i}]_{(\ell,i) \in [L] \times [H]} \in \mathbb{R}^{2dHL}$ and $\mathbf{v} = [\mathbf{w}_Q^{\ell,i}, \mathbf{w}_O^{\ell,i}]_{(\ell,i) \in [L] \times [H]} \in \mathbb{R}^{2dHL}$, and
 133 some smooth model h , which fits under Definition 3.1.

134 4 Incremental learning in networks with diagonal weights

135 Any model f_{NN} with diagonal weights as in Definition 3.1 evolves under the gradient flow (3) as

$$\begin{aligned} \frac{d\mathbf{u}}{dt} &= \mathbf{v} \odot \mathbf{g}(\theta), & \frac{d\mathbf{v}}{dt} &= \mathbf{u} \odot \mathbf{g}(\theta) \quad \text{where} \\ \mathbf{g}(\theta) &= -\mathbb{E}_{\mathbf{x}, \mathbf{y}} [D\ell(\mathbf{y}, h(\mathbf{x}; \mathbf{u} \odot \mathbf{v}))^\top Dh(\mathbf{x}; \mathbf{u} \odot \mathbf{v})^\top]. \end{aligned} \quad (4)$$

136 Here $D\ell(\mathbf{y}, \cdot) \in \mathbb{R}^{1 \times d_{out}}$ is the derivative of ℓ in the second argument and $Dh(\mathbf{x}, \cdot) \in \mathbb{R}^{d_{out} \times p}$ is
 137 the derivative of h in the second argument. We show that if initialization scale of $\theta = (\mathbf{u}, \mathbf{v})$ is
 138 small, then learning proceeds in incremental stages, as given in Algorithm 1, where in each stage the
 139 effective sparsity of \mathbf{u} and \mathbf{v} increases by at most one.

140 4.1 Intuition for incremental learning dynamics

141 We develop an informal intuition for the result. First, we observe a conservation law that simplifies
 142 the dynamics. It can be viewed as the balancedness property for networks with linear activations
 143 Arora et al. (2018); Du et al. (2018), specialized to the case of diagonal layers.

144 **Lemma 4.1** (Conservation law). For any $i \in [p]$ and any time t , we have

$$u_i^2(t) - v_i^2(t) = u_i^2(0) - v_i^2(0). \quad (5)$$

145 *Proof.* This follows from $\frac{d}{dt}(u_i^2 - v_i^2) = u_i v_i g_i(\theta) - u_i v_i g_i(\theta) = 0$. \square

146 This reduces the degrees of freedom and means that we need only keep track of p parameters in total.
 147 Specifically, if we define $w_i(t) := u_i(t) + v_i(t)$, then the vector $\mathbf{w} = \mathbf{u} + \mathbf{v}$ evolves by

$$\frac{d\mathbf{w}}{dt} = \mathbf{w} \odot \mathbf{g}(\theta). \quad (6)$$

148 Using the conservation law (5), one can compute $\theta(t)$ from $\mathbf{w}(t)$, so it remains to analyze the
 149 dynamics of $\mathbf{w}(t)$.

150 **4.1.1 Stage 1 of dynamics**

Stage 1A of dynamics: loss plateau for time $\Theta(\log(1/\alpha))$ At very early times t , we have $\boldsymbol{\theta}(t) \approx \mathbf{0}$ because the weights are initialized to be very small. Thus, we can approximate $\mathbf{g}(\boldsymbol{\theta}(t)) \approx \mathbf{g}(\mathbf{0})$ and so we can solve for the evolution of \mathbf{w} :

$$\mathbf{w}(t) \approx \mathbf{w}(0) \odot e^{\mathbf{g}(\mathbf{0})t}.$$

151 This approximation is valid until one of the entries of $\boldsymbol{\theta}(t)$ reaches constant size, which one can show
152 happens around time $t \approx T_1 \cdot \log(1/\alpha)$ for

$$T_1 = \min_{i \in [p]} 1/|g_i(\mathbf{0})|.$$

153 Until this time, the weights $\boldsymbol{\theta}(t)$ are small, the network remains close to its initialization, and so we
154 observe a loss plateau.

155 **Stage 1B of dynamics: nonlinear dynamics for time $O(1)$** Subsequently, we observe a rapid
156 decrease of the loss and nonlinear dynamics during a $O(1)$ -order time-scale. Indeed, suppose
157 that the dynamics are “non-degenerate” in the sense that there is a unique coordinate i_0 such that
158 $1/|g_{i_0}(\mathbf{0})| = T_1$. Under this assumption, in stage 1A, the weights only grow significantly at
159 coordinate i_0 . So one can show that for any small $\epsilon > 0$, there is a time $t_1(\epsilon) \approx T_1 \cdot \log(1/\alpha)$ such
160 that $u_{i_0}(t_1) \approx \epsilon$, $v_{i_0}(t_1) \approx s\epsilon$ for some sign $s \in \{+1, -1\}$, and $|u_i(t_1)|, |v_i(t_1)| = o_\alpha(1)$ for all
161 $i \neq i_0$.³

162 Because all coordinates except for i_0 are negligibly small after stage 1A, we may perform the
163 following approximation of the dynamics. Zero out the weights at coordinates except for i_0 , and
164 consider the training dynamics starting at $\tilde{\boldsymbol{\theta}} = (\epsilon e_{i_0}, s\epsilon e_{i_0})$. After some constant time, independent
165 of α , these dynamics should approach a stationary point. Furthermore, all coordinates of \mathbf{u} and \mathbf{v}
166 will remain zero except for the i_0 coordinate, so the sparsity of the weights will be preserved. In other
167 words, we should expect there to be a time $\bar{t}_1 = t_1 + O(1) \approx T_1 \cdot \log(1/\alpha)$ such that

$$\boldsymbol{\theta}(\bar{t}_1) \approx (a e_{i_0}, s a e_{i_0}) := \boldsymbol{\theta}^1,$$

168 for some $a \in \mathbb{R}_{>0}$, such that $\boldsymbol{\theta}^1$ is a stationary point of the loss.⁴ This is a good approximation
169 because $\bar{t}_1 - t_1 = O(1)$ is a constant time-scale, so the weights at coordinates except for i_0 remain
170 negligible between times t_1 and \bar{t}_1 . Overall, we have argued that the network approximately reaches
171 stationary point that is 1-sparse, where only the weights at coordinate i_0 are nonzero.

172 **4.1.2 Later stages**

173 We can extend the argument to any number of stages k , where in each stage the weights remain close
174 to constant for time $\Theta(\log(1/\alpha))$ and then rapidly change during time $O(1)$, with the sparsity of the
175 weights increasing by at most one. In order to analyze multiple stages, we must also keep track of the
176 magnitude of the weights on the logarithmic scale because these evolve nonnegligibly throughout
177 training. Inductively on k , suppose that there is some $T_k \in \mathbb{R}$, $\mathbf{b}^k \in \mathbb{R}^p$ and $\boldsymbol{\theta}^k \in \mathbb{R}^{2p}$ and a time
178 $\bar{t}_k \approx T_k \cdot \log(1/\alpha)$ such that

$$\log_\alpha(\mathbf{w}(\bar{t}_k)) \approx \mathbf{b}^k \text{ and } \boldsymbol{\theta}(\bar{t}_k) \approx \boldsymbol{\theta}^k,$$

179 where $\boldsymbol{\theta}^k$ is a stationary point of the loss. We argue for the inductive step that there is $T_{k+1} \in \mathbb{R}$ such
180 that during times $t \in (\bar{t}_k, T_{k+1} \cdot \log(1/\alpha) - \Omega(1))$ the weights remain close to the stationary point
181 from the previous phase, i.e., $\boldsymbol{\theta}(t) \approx \boldsymbol{\theta}^k$. And at a time $\bar{t}_{k+1} \approx T_{k+1} \cdot \log(1/\alpha)$ we have

$$\log_\alpha(\mathbf{w}(\bar{t}_{k+1})) \approx \mathbf{b}^{k+1} \text{ and } \boldsymbol{\theta}(\bar{t}_{k+1}) \approx \boldsymbol{\theta}^{k+1},$$

182 where $\boldsymbol{\theta}^{k+1}$ and \mathbf{b}^{k+1} are defined below, and $\boldsymbol{\theta}^{k+1}$ is a stationary point of the loss whose support
183 has grown by at most one compared to $\boldsymbol{\theta}^k$. The pseudocode for the evolution of \mathbf{b}^k and $\boldsymbol{\theta}^k$ along the
184 stages is given in Algorithm 1, and more details are provided below.

³Without loss of generality, we can ensure that at initialization $\mathbf{u}(0)$ and $\mathbf{u}(0) + \mathbf{v}(0)$ are nonnegative. This implies $\mathbf{u}(t)$ is nonnegative. The fact that u_{i_0} and v_{i_0} are roughly equal in magnitude but might differ in sign is due to the conservation law (5). See Appendix A.3 for details.

⁴The entries of \mathbf{u} and \mathbf{v} are close in magnitude (but may differ in sign) because of the conservation law (5).

185 **Stage $(k + 1)$ A, loss plateau for time $\Theta(\log(1/\alpha))$** At the beginning of stage $k + 1$, the weights
 186 are close to the stationary point θ^k , and so, similarly to stage 1A, linear dynamics are valid.

$$\mathbf{w}(t) \approx \mathbf{w}(\bar{t}_k) \odot e^{\mathbf{g}(\theta^k)(t-\bar{t}_k)}. \quad (7)$$

187 Using the conservation law (5), we derive a “time until active” for each coordinate $i \in [p]$, which
 188 corresponds to the time for the weight at that coordinate to grow from negligible to nonnegligible
 189 magnitude:

$$\Delta_k(i) = \begin{cases} (b_i^k - 1 + \text{sgn}(g_i(\theta^k)))/g_i(\theta^k), & \text{if } g_i(\theta^k) \neq 0 \\ \infty, & \text{if } g_i(\theta^k) = 0 \end{cases} \quad (8)$$

190 The approximation (7) therefore breaks down at a time $t \approx T_{k+1} \cdot \log(1/\alpha)$, where

$$T_{k+1} = T_k + \Delta_k(i_k), \quad i_k = \arg \min_{i \in [p]} \Delta_k(i), \quad (9)$$

191 which corresponds to the first time at the weights at a coordinate grow from negligible to nonnegligible
 192 magnitude. And at times $t \approx T_{k+1} \cdot \log(1/\alpha)$, on the logarithmic scale \mathbf{w} is given by

$$\log_\alpha(\mathbf{w}(t)) \approx \mathbf{b}^{k+1} := \mathbf{b}^k - \mathbf{g}(\theta^k)\Delta_k(i_k), \quad (10)$$

193 **Stage $(k + 1)$ B of dynamics: nonlinear dynamics for time $O(1)$** Subsequently, the weights evolve
 194 nonlinearly during $O(1)$ time. To see this, if we make the non-degeneracy assumption that there
 195 is a unique coordinate i_k such that $\Delta_k(i_k) = \min_i \Delta_k(i)$, then this means that in stage $(k + 1)$ A,
 196 the only coordinate where weights grow from negligible to nonnegligible magnitude is i_k . Roughly
 197 speaking, for any $\epsilon > 0$, there is a time $\underline{t}_{k+1}(\epsilon) \approx T_{k+1} \cdot \log(1/\alpha)$ such that

$$\theta(\underline{t}_{k+1}) \approx \theta^k + (\epsilon e_{i_k}, \text{sgn}(g_{i_k}(\theta^k))\epsilon e_{i_k}),$$

198 where the sign of the weights in coordinate i_k comes from the conservation law (5). At this time,
 199 the weights are approximately the stationary point from stage k , plus a small perturbation. Consider
 200 the dynamics of $\psi^k(t, \epsilon) \in \mathbb{R}^{2p}$ initialized at $\psi^k(0, \epsilon) = \theta^k + (\epsilon e_{i_k}, \text{sgn}(g_{i_k}(\theta^k))\epsilon e_{i_k})$ and evolving
 201 according to the gradient flow $\frac{d\psi^k(t, \epsilon)}{dt} = -\nabla_{\theta} \mathcal{L}(\psi^k)$. These dynamics may be highly nonlinear, so
 202 to control them let us assume that as we take ϵ to be small, they converge to a limiting point θ^{k+1}

$$\lim_{\epsilon \rightarrow 0} \lim_{t \rightarrow \infty} \psi^k(t, \epsilon) = \theta^{k+1}. \quad (11)$$

203 Then we expect that at a time $\bar{t}_{k+1} = \underline{t}_{k+1} + O(1) \approx T_{k+1} \cdot \log(1/\alpha)$, we have $\theta(\bar{t}_{k+1}) \approx \theta^{k+1}$.
 204 This concludes the inductive step.

205 4.2 Formal statement of incremental dynamics

206 We formally state our result. For ease of notation, we write $\theta^k = (\mathbf{u}^k, \mathbf{v}^k)$ and $\mathbf{v}^k = \mathbf{s}^k \odot \mathbf{u}^k$ for
 207 some sign-flip vector $\mathbf{s}^k \in \{+1, -1\}^k$. This form of θ^k can be guaranteed by the conservation law
 208 (5) of the dynamics; see Appendix A. We also denote $\text{supp}(\theta^k) := \text{supp}(\mathbf{u}^k) = \text{supp}(\mathbf{v}^k) \subseteq [p]$.

209 We state our assumptions formally. First, we require that the dynamics be non-degenerate, in the
 210 sense that two coordinates do not become active at the same time. We also place a technical condition
 211 to handle the corner case when a coordinate leaves the support of active coordinates.

Algorithm 1 Incremental learning in networks with diagonal weights

- 1: $\mathbf{b}^0, \theta^0 \leftarrow \mathbf{0} \in \mathbb{R}^p, T_0 \leftarrow 0$
 - 2: **for** stage number $k = 0, 1, 2, \dots$ **do**
 - 3: # (A) Pick new coordinate $i_k \in [p]$ to activate.
 - 4: For each i , define time $\Delta_k(i)$ until active using (8).
 - 5: Pick winning coordinate i_k using (9)
 - 6: Calculate time T_{k+1} using (9) and **break** if ∞
 - 7: Update logarithmic weight approximation \mathbf{b}^{k+1} using (10)
 - 8: # (B) Train activated coordinates to stationarity.
 - 9: $\theta^{k+1} \leftarrow$ limiting dynamics point from (11)
 - 10: **end for**
-

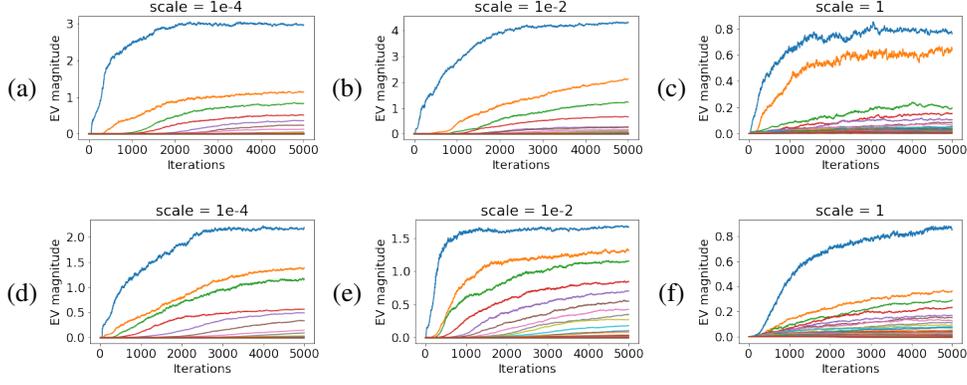


Figure 2: Training a vision transformer on CIFAR-10 using Adam, while varying the initialization scale (unit scale indicates default initialization). Plotted are the evolution of the eigenvalues of $\Delta \mathbf{W}_K \mathbf{W}_Q^\top$ (a) - (c) and $\Delta \mathbf{W}_V \mathbf{W}_O^\top$ (d) - (f) in a random self-attention head in the second layer throughout training. Incremental learning dynamics and a low-rank bias are evident for all scales, albeit more pronounced at smaller initialization scales.

212 **Assumption 4.2** (Nondegeneracy of dynamics in part (A)). The initialization satisfies $u_i(0) \neq v_i(0)$
 213 for all i . For stage k , either $T_k = \infty$ or there is a unique minimizer i_k to $\min_i \Delta_k(i_k)$ in (9). Finally,
 214 for all $i \in \text{supp}(\boldsymbol{\theta}^{k-1}) \setminus \text{supp}(\boldsymbol{\theta}^k)$ we have $g_i(\boldsymbol{\theta}^k) \neq 0$.

215 Next, we require that very small perturbations of the coordinates outside of $\text{supp}(\boldsymbol{\theta}^k)$ do not change
 216 the dynamics. For this, it suffices that $\boldsymbol{\theta}^k$ be a strict local minimum.

217 **Assumption 4.3** (Stationary points are strict local minima). For stage k , there exist $\delta_k > 0$ and
 218 $c_k > 0$ such that for $\tilde{\mathbf{u}} \in B(\mathbf{u}^k, \delta)$ supported on $\text{supp}(\mathbf{u}^k)$, we have

$$\mathcal{L}(\tilde{\mathbf{u}}, \mathbf{s}^k \odot \tilde{\mathbf{u}}) \geq c_k \|\mathbf{u}^k - \tilde{\mathbf{u}}\|^2$$

219 Finally, we require a robust version of the assumption (11), asking for convergence to a neighborhood
 220 of $\boldsymbol{\theta}^{k+1}$ even when the initialization is slightly noisy.

221 **Assumption 4.4** (Noise-robustness of dynamics in part (B)). For any stage k with $T_{k+1} < \infty$ and any
 222 $\epsilon > 0$, there are $\delta > 0$ and $\tau : \mathbb{R}_{>0} \rightarrow \mathbb{R}$ such that the following holds. For any $\tilde{\mathbf{u}} \in B(\mathbf{u}^k, \delta) \cap \mathbb{R}_{\geq 0}^p$
 223 supported on $\text{supp}(\tilde{\mathbf{u}}) \subseteq \text{supp}(\mathbf{u}^k) \cup \{i_k\}$, there exists a unique solution $\boldsymbol{\psi} : [0, \infty) \rightarrow \mathbb{R}^p$ of the
 224 gradient flow $\frac{d\boldsymbol{\psi}}{dt} = -\nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\psi})$ initialized at $\boldsymbol{\psi}(0) = (\tilde{\mathbf{u}}, \mathbf{s}^{k+1} \odot \tilde{\mathbf{u}})$, and at times $t \geq \tau(\tilde{\boldsymbol{\psi}}_{i_k})$,

$$\|\boldsymbol{\psi}(t) - \boldsymbol{\theta}^{k+1}\| < \epsilon.$$

225 These assumptions are validated experimentally in Appendix C. Using them, we prove that incremen-
 226 tal learning Algorithm 1 tracks the gradient flow dynamics if the initialization scale is small.

227 **Theorem 4.5** (Incremental dynamics with untied weights). For any stage k and time $t \in (T_k, T_{k+1})$
 228 the following holds under Assumptions 4.2 4.3 and 4.4. There is $\alpha_0(t) > 0$ such that for all $\alpha < \alpha_0$,
 229 there exists a unique solution $\boldsymbol{\theta} : [0, t \log(1/\alpha)] \rightarrow \mathbb{R}^p$ to the gradient flow (3) and

$$\lim_{\alpha \rightarrow 0} \boldsymbol{\theta}(t \cdot \log(1/\alpha)) \rightarrow \boldsymbol{\theta}^k,$$

230 and at each stage the sparsity increases by at most one: $\text{supp}(\boldsymbol{\theta}^{k+1}) \setminus \text{supp}(\boldsymbol{\theta}^k) \subseteq \{i_k\}$.

231 **Example 4.6** (Application: Incremental learning in diagonal transformer). In Example 3.2, we
 232 showed that a diagonal transformer falls under Theorem 4.5. As a corollary, the gradient flow on a
 233 transformer with small initialization will learn in stages, where in each stage there will be at most one
 234 head $i \in [H]$ on one layer $\ell \in [L]$ such that either the rank of $\mathbf{W}_K^{\ell,i} (\mathbf{W}_Q^{\ell,i})^\top = \text{diag}(\mathbf{w}_K^{\ell,i}) \text{diag}(\mathbf{w}_Q^{\ell,i})$
 235 or the rank of $\mathbf{W}_V^{\ell,i} (\mathbf{W}_O^{\ell,i})^\top = \text{diag}(\mathbf{w}_V^{\ell,i}) \text{diag}(\mathbf{w}_O^{\ell,i})$ increases by at most one.

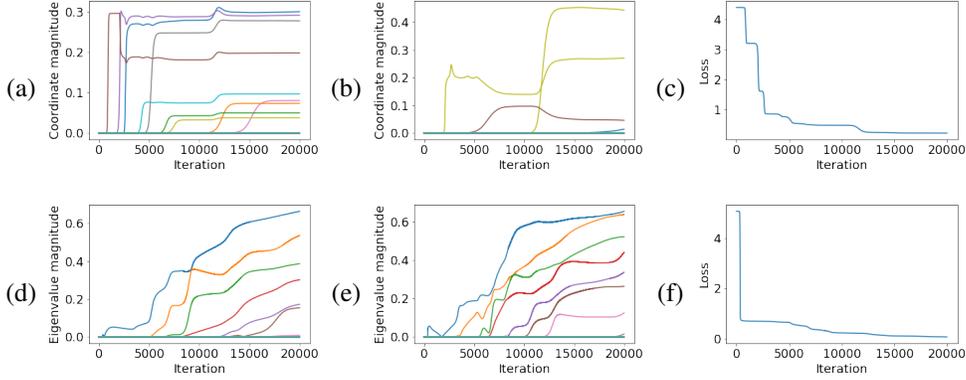


Figure 3: A network containing a single self-attention layer with diagonal (a) - (c) and full (d) - (f) weight matrices, trained with gradient descent in the incremental learning regime. (a) The diagonal entries of $\mathbf{W}_V \mathbf{W}_O^\top$ and (d) the singular values of $\mathbf{W}_V \mathbf{W}_O^\top$ are learned incrementally. (b) The diagonal entries of $\mathbf{W}_K \mathbf{W}_Q^\top$ and (e) the singular values of $\mathbf{W}_K \mathbf{W}_Q^\top$ are learned incrementally. (c), (f) The loss curves show stagewise plateaus and sharp decreases.

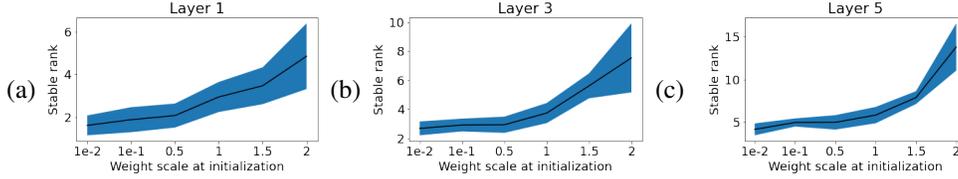


Figure 4: Stable rank of $\Delta \mathbf{W}_K \mathbf{W}_Q^\top$ per initialization scale (Unit scale refers to the default initialization) in different self-attention heads post-training, at layers 1, 3, 5. At each layer, the stable rank mean and standard deviation are computed across 8 heads per layer, for each initialization scale. All models were trained on CIFAR-10 using the Adam optimizer. Smaller initialization scales lead to lower-rank attention heads. Analogous plots for $\Delta \mathbf{W}_V \mathbf{W}_O^\top$ are in the appendix.

236 5 Experimental results

237 We experimentally support our theoretical findings in a series of experiments: first on a toy model
 238 given by Equation (1), followed by experiments on a vision transformer on the CIFAR datasets. We
 239 defer additional experimental details and results to the appendix.

240 **Toy models** We consider a toy model comprised of one self-attention layer with a single head as in
 241 (1), with either diagonal or full weight matrices. We initialize $\mathbf{W}_K, \mathbf{W}_Q, \mathbf{W}_V, \mathbf{W}_O$ using Gaussian
 242 initialization with a small standard deviation, and train the model using GD on a regression task with
 243 50-dimensional random Gaussian token inputs and targets from a teacher model. During training,
 244 we track the diagonal entries of $\mathbf{W}_K \mathbf{W}_Q^\top$ and $\mathbf{W}_V \mathbf{W}_O^\top$ in the diagonal case, and the singular values
 245 of $\mathbf{W}_K \mathbf{W}_Q^\top$ and $\mathbf{W}_V \mathbf{W}_O^\top$ in the full weights case. Our results are summarized in Figure 3. For the
 246 diagonal model, as predicted, diagonal components are learned incrementally, resulting in progressive
 247 increase in the rank; in Appendix C we run additional experiments to verify that the assumptions of
 248 Theorem 4.5 indeed hold. For the full-weights model, we also observe incremental learning with
 249 progressively-increasing rank, even though this setting falls beyond our theory.

250 **Vision transformers** We next run experiments that go well beyond our toy model to test the
 251 extent to which incremental learning with a low-rank bias exists in popular models used in practice.
 252 We conduct experiments with vision transformers (ViT) Dosovitskiy et al. (2020) trained on the
 253 CIFAR-10/100 and ImageNet datasets. We use a ViT of depth 6, with 8 self-attention heads per layer
 254 (with layer normalization). We use an embedding and MLP dimension of $d_{\text{emb}} = 512$, and a head
 255 dimension of $d_h = 128$ (i.e. $\mathbf{W}_K, \mathbf{W}_Q, \mathbf{W}_V, \mathbf{W}_O \in \mathbb{R}^{d_{\text{emb}} \times d_h}$). We train the transformer using Adam

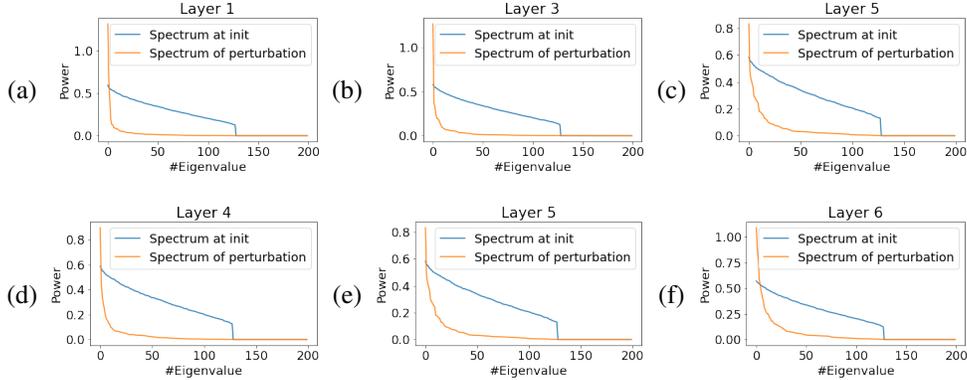


Figure 5: Spectrum of the weight perturbation $\Delta\mathbf{W}_K\mathbf{W}_Q^\top$ vs. initialization in a vision transformer trained on CIFAR-10, using Adam and default initialization scale, in random self-attention heads in different layers. The learned perturbation exhibits extreme low-rank bias post-training even in default initialization scales. Analogous plots for $\Delta\mathbf{W}_V\mathbf{W}_O^\top$ are in the appendix.

256 on the CIFAR-10/100 and ImageNet classification tasks with cross-entropy loss. We train all layers
 257 (including the feedforward layers) while varying the initialization scale of all layers by multiplying
 258 their initial values by a scale factor (we fix the scale of the initial token mapper). To illustrate
 259 the effect of training on weights with a non-vanishing initialization scale, we plot the spectrum
 260 of the difference $\Delta\mathbf{W}_K\mathbf{W}_Q^\top$ and $\Delta\mathbf{W}_V\mathbf{W}_O^\top$ between the weights post-training, and their initial
 261 values. Figure 2 shows the evolution of the principal components of $\Delta\mathbf{W}_K\mathbf{W}_Q^\top$ and $\Delta\mathbf{W}_V\mathbf{W}_O^\top$ for
 262 a randomly-chosen self-attention head and layer throughout training, exhibiting incremental learning
 263 dynamics and a low-rank bias. Note that incremental learning and low-rank bias are increasingly
 264 evident with smaller initialization scales, as further demonstrated in Figure 4. Finally, we plot the
 265 spectrum of $\Delta\mathbf{W}_K\mathbf{W}_Q^\top$ against that of its initialized state in Figure 5 for different self-attention heads,
 266 illustrating that the weight perturbation learned during the training process is extremely low-rank
 267 when compared to the initial spectrum. All figures in this section are given for models trained on
 268 CIFAR-10. In the appendix we conduct further experiments on CIFAR-100 and ImageNet, as well as
 269 different model sizes for completeness, and these show similar trends. Further experimental details
 270 and results are provided in the appendix.

271 6 Discussion

272 We have identified incremental learning dynamics in transformers, proved them rigorously in a
 273 simplified setting, and shown them experimentally in networks trained with practical hyperparameters.

274 **Limitations** There are clear limitations to our theory: the diagonal weights and small initialization
 275 assumptions. More subtly, the theory does not apply to losses with exponential-like tails because the
 276 weights may not converge to a finite value and so Assumption 4.3 is not met (this could possibly be
 277 addressed by adding regularization). Also, the architecture must be smooth, which precludes ReLUs –
 278 but allows for smoothed ReLUs such as the GeLUs used in ViT (Dosovitskiy et al., 2020). Finally,
 279 the theory is for training with gradient flow, while other optimizers such as Adam are used in practice
 280 instead (Kingma & Ba, 2014). Nevertheless, our experiments on ViTs indicate that the incremental
 281 learning dynamics occurs even when training with Adam.

282 **Future directions** A promising direction of future research is to examine the connection between
 283 our results on incremental dynamics and the LoRA method (Hu et al., 2021), with the goal of
 284 explaining and improving on this algorithm; see also the discussion in Section 1.1. Another interesting
 285 avenue is to develop a theoretical understanding of the implicit bias in function space of transformers
 286 whose weights are a low-rank perturbation of randomly initialized weights.

287 **References**

- 288 Abbe, E., Boix-Adsera, E., and Misiakiewicz, T. The merged-staircase property: a necessary and
289 nearly sufficient condition for SGD learning of sparse functions on two-layer neural networks,
290 COLT, 2022.
- 291 Abbe, E., Boix-Adsera, E., and Misiakiewicz, T. Sgd learning on neural networks: leap complexity
292 and saddle-to-saddle dynamics. *arXiv preprint arXiv:2302.11055*, 2023.
- 293 Aghajanyan, A., Zettlemoyer, L., and Gupta, S. Intrinsic dimensionality explains the effectiveness of
294 language model fine-tuning. *arXiv preprint arXiv:2012.13255*, 2020.
- 295 Arora, S., Cohen, N., and Hazan, E. On the optimization of deep networks: Implicit acceleration by
296 overparameterization. In *International Conference on Machine Learning*, pp. 244–253. PMLR,
297 2018.
- 298 Arora, S., Cohen, N., Hu, W., and Luo, Y. Implicit regularization in deep matrix factorization.
299 *Advances in Neural Information Processing Systems*, 32, 2019.
- 300 Berthier, R. Incremental learning in diagonal linear networks. *arXiv preprint arXiv:2208.14673*,
301 2022.
- 302 Bietti, A., Bruna, J., Sanford, C., and Song, M. J. Learning single-index models with shallow neural
303 networks. *arXiv preprint arXiv:2210.15651*, 2022.
- 304 Boursier, E., Pillaud-Vivien, L., and Flammarion, N. Gradient flow dynamics of shallow relu networks
305 for square loss and orthogonal inputs. *arXiv preprint arXiv:2206.00939*, 2022.
- 306 Chizat, L., Oyallon, E., and Bach, F. R. On lazy training in differentiable programming. In *Neural
307 Information Processing Systems*, 2018.
- 308 Damian, A., Lee, J., and Soltanolkotabi, M. Neural networks can learn representations with gradient
309 descent. In *Conference on Learning Theory*, pp. 5413–5452. PMLR, 2022.
- 310 Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. Bert: Pre-training of deep bidirectional
311 transformers for language understanding. *ArXiv*, abs/1810.04805, 2019.
- 312 Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani,
313 M., Minderer, M., Heigold, G., Gelly, S., Uszkoreit, J., and Houlsby, N. An image is worth 16x16
314 words: Transformers for image recognition at scale. *ArXiv*, abs/2010.11929, 2020.
- 315 Du, S. S., Hu, W., and Lee, J. D. Algorithmic regularization in learning deep homogeneous models:
316 Layers are automatically balanced. *Advances in Neural Information Processing Systems*, 31, 2018.
- 317 Frei, S., Vardi, G., Bartlett, P. L., Srebro, N., and Hu, W. Implicit bias in leaky relu networks trained
318 on high-dimensional data. *arXiv preprint arXiv:2210.07082*, 2022.
- 319 Geiger, M., Spigler, S., Jacot, A., and Wyart, M. Disentangling feature and lazy learning in deep
320 neural networks: an empirical study. *ArXiv*, abs/1906.08034, 2019.
- 321 Ghorbani, B., Mei, S., Misiakiewicz, T., and Montanari, A. Limitations of lazy training of two-layers
322 neural network. *Advances in Neural Information Processing Systems*, 32, 2019.
- 323 Gidel, G., Bach, F., and Lacoste-Julien, S. Implicit regularization of discrete gradient dynamics in
324 linear neural networks. *Advances in Neural Information Processing Systems*, 32, 2019.
- 325 Gissin, D., Shalev-Shwartz, S., and Daniely, A. The implicit bias of depth: How incremental learning
326 drives generalization. *arXiv preprint arXiv:1909.12051*, 2019.
- 327 Hu, E. J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., Wang, L., and Chen, W. Lora:
328 Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021.
- 329 Jacot, A., Gabriel, F., and Hongler, C. Neural tangent kernel: Convergence and generalization in
330 neural networks. *Advances in neural information processing systems*, 31, 2018.

- 331 Jacot, A., Ged, F. G., Simsek, B., Hongler, C., and Gabriel, F. Saddle-to-saddle dynamics in deep
332 linear networks: Small initialization training, symmetry, and sparsity. 2021.
- 333 Jin, J., Li, Z., Lyu, K., Du, S. S., and Lee, J. D. Understanding incremental learning of gradient
334 descent: A fine-grained analysis of matrix sensing. *arXiv preprint arXiv:2301.11500*, 2023.
- 335 Kim, D. and Chung, H. W. Rank-1 matrix completion with gradient descent and small random
336 initialization. *ArXiv*, abs/2212.09396, 2022.
- 337 Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization. *arXiv preprint*
338 *arXiv:1412.6980*, 2014.
- 339 Li, C., Farkhoor, H., Liu, R., and Yosinski, J. Measuring the intrinsic dimension of objective
340 landscapes. *arXiv preprint arXiv:1804.08838*, 2018.
- 341 Li, Z., Luo, Y., and Lyu, K. Towards resolving the implicit bias of gradient descent for matrix
342 factorization: Greedy low-rank learning. *ArXiv*, abs/2012.09839, 2020.
- 343 Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., Levy, O., Lewis, M., Zettlemoyer, L., and
344 Stoyanov, V. Roberta: A robustly optimized bert pretraining approach. *ArXiv*, abs/1907.11692,
345 2019.
- 346 Malach, E., Kamath, P., Abbe, E., and Srebro, N. Quantifying the benefit of using differentiable
347 learning over tangent kernels. In Meila, M. and Zhang, T. (eds.), *Proceedings of the 38th Inter-*
348 *national Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning*
349 *Research*, pp. 7379–7389. PMLR, 18–24 Jul 2021. URL <https://proceedings.mlr.press/v139/malach21a.html>.
350
- 351 Milanese, P., Kadri, H., Ayache, S., and Artières, T. Implicit regularization in deep tensor factorization.
352 *2021 International Joint Conference on Neural Networks (IJCNN)*, pp. 1–8, 2021.
- 353 Moroshko, E., Gunasekar, S., Woodworth, B. E., Lee, J., Srebro, N., and Soudry, D. Implicit bias in
354 deep linear classification: Initialization scale vs training accuracy. *ArXiv*, abs/2007.06738, 2020.
- 355 Mousavi-Hosseini, A., Park, S., Girotti, M., Mitliagkas, I., and Erdogdu, M. A. Neural networks
356 efficiently learn low-dimensional representations with sgd. *arXiv preprint arXiv:2209.14863*, 2022.
- 357 Patel, D. and Ahmad, A. Google “we have no moat, and neither does openai”, May 2023. URL
358 <https://www.semianalysis.com/p/google-we-have-no-moat-and-neither>.
- 359 Stöger, D. and Soltanolkotabi, M. Small random initialization is akin to spectral learning: Opti-
360 mization and generalization guarantees for overparameterized low-rank matrix reconstruction. In
361 *Neural Information Processing Systems*, 2021.
- 362 Taori, R., Gulrajani, I., Zhang, T., Dubois, Y., Li, X., Guestrin, C., Liang, P., and Hashimoto,
363 T. B. Alpaca: A strong, replicable instruction-following model. *Stanford Center for Research on*
364 *Foundation Models*. <https://crfm.stanford.edu/2023/03/13/alpaca.html>, 2023.
- 365 Timor, N., Vardi, G., and Shamir, O. Implicit regularization towards rank minimization in relu
366 networks. In *International Conference on Algorithmic Learning Theory*, pp. 1429–1459. PMLR,
367 2023.
- 368 Vaswani, A., Shazeer, N. M., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L., and
369 Polosukhin, I. Attention is all you need. *ArXiv*, abs/1706.03762, 2017.
- 370 Woodworth, B. E., Gunasekar, S., Lee, J., Moroshko, E., Savarese, P. H. P., Golan, I., Soudry, D., and
371 Srebro, N. Kernel and rich regimes in overparametrized models. *ArXiv*, abs/2002.09277, 2019.
- 372 Yu, H. and Wu, J. Compressing transformers: Features are low-rank, but weights are not! 2023.

373 **A Proof for dynamics of networks with diagonal parametrization**
 374 **(Theorem 4.5)**

375 **A.1 Assumptions**

376 Recall we have defined $\theta^0, \dots, \theta^k, \dots \in \mathbb{R}^{2p}$ as the sequence of weights such that $\theta^0 = \mathbf{0}$ and θ^{k+1}
 377 is defined inductively as follows. Consider the dynamics of $\psi^k(t, \epsilon) \in \mathbb{R}^{2p}$ initialized at $\psi^k(0, \epsilon) =$
 378 $\theta^k + (\epsilon e_{i_k}, \text{sgn}(g_i(\theta^k)) \epsilon e_{i_k})$ and evolving according to the gradient flow $\frac{d\psi^k(t, \epsilon)}{dt} = -\nabla_{\theta} \mathcal{L}(\psi^k)$.
 379 We assume that there is a limiting point θ^{k+1} of these dynamics as ϵ is taken small and the time is
 380 taken large:

$$\lim_{\epsilon \rightarrow 0} \lim_{t \rightarrow \infty} \psi^k(t, \epsilon) = \theta^{k+1}.$$

381 Under the above assumption that this sequence $\theta^0, \dots, \theta^k, \dots$ is well-defined, we can derive a useful
 382 property of it for free. Namely, the conservation law (5) implies that $\mathbf{u} \odot \mathbf{u} - \mathbf{v} \odot \mathbf{v}$ is preserved. It
 383 follows that for each k we have that $\theta^k = (\mathbf{u}^k, \mathbf{v}^k)$ satisfies $|\mathbf{u}^k| = |\mathbf{v}^k|$ entrywise. In other words,
 384 there is $\mathbf{s}^k \in \{+1, -1\}^p$ satisfying

$$\theta^k = (\mathbf{u}^k, \mathbf{s}^k \odot \mathbf{u}^k) \in \mathbb{R}^{2p}.$$

385 We also abuse notation and write $\text{supp}(\theta^k) := \text{supp}(\mathbf{u}^k) \subseteq [p]$, since the support of θ^k on the first p
 386 coordinates matches its support on the last p coordinates.

387 Having fixed this notation, we now recall the main assumptions of the theorem.

388 **Assumption A.1** (Nondegeneracy of dynamics in part (A); Assumption 4.2). The initialization
 389 satisfies $u_i(0) \neq v_i(0)$ for all i . For stage k , either $T_{k+1} = \infty$ or there is a unique minimizer i_k to
 390 $\min_i \Delta_k(i_k)$ in (9). Finally, for all $i \in \text{supp}(\theta^{k-1}) \setminus \text{supp}(\theta^k)$ we have $g_i(\theta^k) \neq 0$.

391 **Assumption A.2** (Stationary points are strict local minima; Assumption 4.3). For stage k , there exist
 392 $\delta_k > 0$ and $c_k > 0$ such that for $\tilde{\mathbf{u}} \in B(\mathbf{u}^k, \delta)$ supported on $\text{supp}(\mathbf{u}^k)$, we have

$$\mathcal{L}(\tilde{\mathbf{u}}, \mathbf{s}^k \odot \tilde{\mathbf{u}}) \geq c_k \|\mathbf{u}^k - \tilde{\mathbf{u}}\|^2.$$

393 **Assumption A.3** (Noise-robustness of dynamics in part (B); Assumption 4.4). For stage k , either
 394 $T_{k+1} = \infty$ or the following holds. For any $\epsilon > 0$, there are $\delta > 0$ and $\tau : \mathbb{R}_{>0} \rightarrow \mathbb{R}$ such that
 395 the following holds. For any $\tilde{\mathbf{u}} \in B(\mathbf{u}^k, \delta) \cap \mathbb{R}_{\geq 0}^p$ supported on $\text{supp}(\tilde{\mathbf{u}}) \subseteq \text{supp}(\mathbf{u}^k) \cup \{i_k\}$,
 396 there exists a unique solution $\psi : [0, \infty) \rightarrow \mathbb{R}^p$ of the gradient flow $\frac{d\psi}{dt} = -\nabla_{\theta} \mathcal{L}(\psi)$ initialized at
 397 $\psi(0) = (\tilde{\mathbf{u}}, \mathbf{s}^{k+1} \odot \tilde{\mathbf{u}})$, and at times $t \geq \tau(\tilde{u}_{i_k})$,

$$\|\psi(t) - \theta^{k+1}\| < \epsilon.$$

398 **A.2 Rescaling time for notational convenience**

399 For ease of notation, we rescale time

$$\begin{aligned} \mathbf{u}_{\alpha}(0) &= \alpha \mathbf{u}(0), & \mathbf{v}_{\alpha}(0) &= \alpha \mathbf{v}(0) \\ \frac{d\mathbf{u}_{\alpha}}{dt} &= \log(1/\alpha) \mathbf{v}_{\alpha} \odot \mathbf{g}(\mathbf{u}_{\alpha}, \mathbf{v}_{\alpha}), & \frac{d\mathbf{v}_{\alpha}}{dt} &= \log(1/\alpha) \mathbf{u}_{\alpha} \odot \mathbf{g}(\mathbf{u}_{\alpha}, \mathbf{v}_{\alpha}). \end{aligned} \quad (12)$$

400 We also define

$$\theta_{\alpha}(t) = (\mathbf{u}_{\alpha}(t), \mathbf{v}_{\alpha}(t)) \in \mathbb{R}^{2p}.$$

401 Because of this time-rescaling, we equivalently state Theorem 4.5 as:

402 **Theorem A.4** (Restatement of Theorem 4.5). *Let $K \in \mathbb{Z}_{\geq 0}$ be such that Assumptions 4.2–4.3 hold*
 403 *for all $k \leq K$ and Assumption 4.4 holds for all $k < K$. Then for any $k \leq K$ and time $t \in (T_k, T_{k+1})$*
 404 *the following holds. There is $\alpha_0(t) > 0$ such that for all $\alpha < \alpha_0$, there exists a unique solution*
 405 *$\theta_{\alpha} : [0, t] \rightarrow \mathbb{R}^{2p}$ to the gradient flow (12) and*

$$\lim_{\alpha \rightarrow 0} \theta_{\alpha}(t) \rightarrow \theta^k,$$

406 *where at each stage $|\text{supp}(\mathbf{u}^k) \setminus \text{supp}(\mathbf{u}^{k-1})| \leq 1$.*

For shorthand, we also write

$$S_k = \text{supp}(\mathbf{u}^k) \text{ and } S_k^c = [p] \setminus \text{supp}(\mathbf{u}^k).$$

407 **A.3 Simplifying problem without loss of generality**

408 For each coordinate $i \in [p]$ we have $|u_{\alpha,i}(0)| \neq |v_{\alpha,i}(0)|$ by the non-degeneracy Assumption 4.2.
 409 So we can assume $|u_{\alpha,i}(0)| > |v_{\alpha,i}(0)|$ without loss of generality. Furthermore, we can assume the
 410 entrywise inequality

$$\mathbf{u}_\alpha(0) > 0$$

411 by otherwise training weights $\tilde{\mathbf{u}}_\alpha(t), \tilde{\mathbf{v}}_\alpha(t)$ initialized at $\tilde{\mathbf{u}}_\alpha(0) = \text{sgn}(\mathbf{u}_\alpha(0))\mathbf{u}_\alpha(0)$ and $\tilde{\mathbf{v}}_\alpha(0) =$
 412 $\text{sgn}(\mathbf{v}_\alpha(0))\mathbf{v}_\alpha(0)$, as $\tilde{\mathbf{u}}_\alpha(t) \odot \tilde{\mathbf{v}}_\alpha(t) = \mathbf{u}_\alpha(t) \odot \mathbf{v}_\alpha(t)$ at all times.

413 Since $u_{\alpha,i}^2(t) - v_{\alpha,i}^2(t) = u_{\alpha,i}^2(0) - v_{\alpha,i}^2(0)$ by the conservation law (5), it holds that $|u_{\alpha,i}(t)| >$
 414 $|v_{\alpha,i}(t)|$ throughout. So by continuity

$$\mathbf{u}_\alpha(t) > 0$$

415 throughout training.

416 **A.4 Tracking the sum of the weights**

417 We define

$$\mathbf{w}_\alpha(t) = \mathbf{u}_\alpha(t) + \mathbf{v}_\alpha(t).$$

418 The reason for this definition is that during training we have

$$\frac{d\mathbf{w}_\alpha}{dt} = \log(1/\alpha)\mathbf{w}_\alpha \odot \mathbf{g}(\boldsymbol{\theta}_\alpha), \quad (13)$$

419 Notice that since that we have assumed $u_{\alpha,i}(0) > |v_{\alpha,i}(0)|$ for each $i \in [p]$ we have $\mathbf{w}_\alpha(0) > 0$
 420 entrywise. So, by (13) for all $t > 0$,

$$\mathbf{w}_\alpha(t) > 0.$$

421 It suffices to track $\mathbf{w}_\alpha(t)$ because we can relate the log-scale magnitude of $\mathbf{w}_\alpha(t)$ to the magnitudes
 422 of the corresponding coordinates in $\mathbf{u}_\alpha(t)$ and $\mathbf{v}_\alpha(t)$ – see technical Lemmas B.1 B.2 and B.3.

423 **A.5 Claimed invariants in proof of Theorem A.4**

424 In order to prove Theorem A.4, we consider any gradient flow $\boldsymbol{\theta}_\alpha : [0, T^*] \rightarrow \mathbb{R}^p$ solving (12) where
 425 $T^* \in (T_K, T_{K+1})$. For now, we focus only on proving properties of this gradient flow, and defer its
 426 existence and uniqueness to Section A.8.

427 We show the following invariants inductively on the stage k . For any $\epsilon > 0$, any stage $k \leq K$, there
 428 is $\alpha_k := \alpha_k(\epsilon) > 0$ such that for all $\alpha < \alpha_k$ the following holds. There are times $\bar{t}_k := \bar{t}_k(\alpha, \epsilon)$ and
 429 $\underline{t}_{k+1} := \underline{t}_{k+1}(\alpha, \epsilon)$, such that

$$\bar{t}_k \in [T_k - \epsilon, T_k + \epsilon], \quad (14)$$

$$\underline{t}_{k+1} \in \begin{cases} [T_{k+1} - \epsilon, T_{k+1} + \epsilon], & \text{if } T_{k+1} < \infty \\ \{T^*\}, & \text{if } T_{k+1} = \infty \end{cases}. \quad (15)$$

430 and the weights approximate the greedy limit for all times $t \in [\bar{t}_k, \underline{t}_{k+1}]$

$$\|\boldsymbol{\theta}_\alpha(t) - \boldsymbol{\theta}^k\| < \epsilon, \quad (16)$$

431 and the weights at times \bar{t}_k and \underline{t}_{k+1} are correctly estimated by the incremental learning dynamics on
 432 the logarithmic-scale

$$\|\log_\alpha(\mathbf{w}_\alpha(\bar{t}_k)) - \mathbf{b}^k\| < \epsilon \quad (17)$$

433 and if $T_{k+1} < \infty$ then

$$\|\log_\alpha(\mathbf{w}_\alpha(\underline{t}_{k+1})) - \mathbf{b}^{k+1}\| < \epsilon. \quad (18)$$

434 *Base case $k = 0$:* Take $\bar{t}_0(\alpha, \epsilon) = 0$. Then statement (14) holds since $T_0 = 0$. Notice that as $\alpha \rightarrow 0$
 435 we have that $\mathbf{u}_\alpha(0), \mathbf{v}_\alpha(0) \rightarrow \mathbf{0} = \mathbf{u}^0$, and also $\log_\alpha \mathbf{w}_\alpha(0) \rightarrow \mathbf{1} = \mathbf{b}^0$. So statement (17) follows if
 436 we take α_0 small enough. In Section A.6 we show how to construct time \underline{t}_1 such that (16) and (18)
 437 hold.

438 *Inductive step:* Suppose that (14), (16), (17) and (18) hold for some iteration $k < K$. We prove them
 439 for iteration $k + 1$. In Section A.7 we construct time \bar{t}_k . In Section A.6 we construct time \underline{t}_{k+1} .

440 **A.6 Dynamics from time \bar{t}_k to time \underline{t}_{k+1} (Linear dynamics for $O(\log(1/\alpha))$ unrescaled time)**

441 Let $k \leq K$, and suppose that we know that for any $\bar{\epsilon}_k > 0$, there is $\bar{\alpha}_k(\bar{\epsilon}_k) > 0$ such that for all
 442 $0 < \alpha < \bar{\alpha}_k$, there is a time $\bar{t}_k = \bar{t}_k(\alpha, \bar{\epsilon}_k)$ satisfying

$$\begin{aligned} |T_k - \bar{t}_k| &< \bar{\epsilon}_k \\ \|\boldsymbol{\theta}_\alpha(\bar{t}_k) - \boldsymbol{\theta}^k\| &< \bar{\epsilon}_k \\ \|\log_\alpha(\mathbf{w}_\alpha(\bar{t}_k)) - \mathbf{b}^k\| &< \bar{\epsilon}_k. \end{aligned}$$

443 **A.6.1 Analysis in case where $T_{k+1} < \infty$**

444 Consider first the case where $T_{k+1} < \infty$. We show that, for any $\underline{\epsilon}_{k+1} > 0$, there is $\rho_{k+1}(\underline{\epsilon}_{k+1}) > 0$
 445 such that for all $0 < \rho < \rho_{k+1}(\underline{\epsilon}_{k+1})$ there is $\underline{\alpha}_{k+1}(\rho, \underline{\epsilon}_{k+1}) > 0$ such that for all $\alpha < \underline{\alpha}_{k+1}$, there
 446 is a time $\underline{t}_{k+1} = \underline{t}_{k+1}(\alpha, \rho, \underline{\epsilon}_{k+1})$ satisfying

$$|T_{k+1} - \underline{t}_{k+1}| < \underline{\epsilon}_{k+1} \tag{19}$$

$$\|\boldsymbol{\theta}_\alpha(t) - \boldsymbol{\theta}^k\| < \underline{\epsilon}_{k+1} \text{ for all } t \in [\bar{t}_k, \underline{t}_{k+1}] \tag{20}$$

$$\|\log_\alpha(\mathbf{w}_\alpha(\underline{t}_{k+1})) - \mathbf{b}^{k+1}\| < \underline{\epsilon}_{k+1} \tag{21}$$

$$u_{\alpha, i_k}(\underline{t}_{k+1}) \in [\rho, 3\rho], \tag{22}$$

$$\text{sgn}(v_{\alpha, i_k}(\underline{t}_{k+1})) = s_{i_k}^{k+1}. \tag{23}$$

447 For any ρ, α , let $\bar{\epsilon}_k = \rho \underline{\epsilon}_{k+1} / (4\rho)$ and choose $\bar{t}_k = \bar{t}_k(\alpha, \bar{\epsilon}_k)$. Then define

$$\begin{aligned} \underline{t}_{k+1} &= \underline{t}_{k+1}(\alpha, \rho, \underline{\epsilon}_{k+1}) \\ &= \inf\{t \in [\bar{t}_k, \infty) : \|\mathbf{u}_{\alpha, S_k^c}(t) - \mathbf{u}_{\alpha, S_k^c}(\bar{t}_k)\| + \|\mathbf{v}_{\alpha, S_k^c}(t) - \mathbf{v}_{\alpha, S_k^c}(\bar{t}_k)\| > 4\rho\}. \end{aligned} \tag{24}$$

448 Now we show that the weights $\boldsymbol{\theta}_\alpha(t)$ cannot move much from time \bar{t}_k to \underline{t}_{k+1} . The argument uses the
 449 local Lipschitzness of the loss \mathcal{L} (from technical Lemma B.7), and the strictness of $\boldsymbol{\theta}^k$ as a stationary
 450 point (from Assumption 4.3).

451 **Lemma A.5** (Stability of active variables during part (A) of dynamics). *There is ρ_{k+1} small enough
 452 and $\underline{\alpha}_{k+1}(\rho)$ small enough depending on ρ , such that for all $\rho < \rho_{k+1}$ and $\alpha < \underline{\alpha}_{k+1}$ and all
 453 $t \in [\bar{t}_k, \underline{t}_{k+1})$,*

$$\|\boldsymbol{\theta}_\alpha(t) - \boldsymbol{\theta}^k\| < \rho' := \max(24\rho, 18\sqrt{\rho K_{R_k}/c_k}). \tag{25}$$

454 where c_k is the strict-minimum constant from Assumption 4.3 and K_{R_k} is the Lipschitzness constant
 455 from Lemma B.7 for the ball of radius $R_k = \|\boldsymbol{\theta}^k\| + 1$.

Proof. Assume by contradiction that (25) is violated at some time $t < \underline{t}_{k+1}$. Let us choose the first
 such time

$$t^* = \inf\{t \in [\bar{t}_k, \underline{t}_{k+1}) : \|\mathbf{u}_\alpha(t^*) - \mathbf{u}^k\| + \|\mathbf{v}_\alpha(t^*) - \mathbf{s}^k \odot \mathbf{u}^k\| \geq \rho'\}.$$

456 Define $\tilde{\boldsymbol{\theta}} = (\tilde{\mathbf{u}}, \tilde{\mathbf{v}})$ by

$$\tilde{u}_i = \begin{cases} u_{\alpha, i}(t^*), & i \in S_k \\ 0, & i \notin S_k \end{cases} \quad \text{and} \quad \tilde{v}_i = \begin{cases} v_{\alpha, i}(t^*), & i \in S_k \\ 0, & i \notin S_k \end{cases}.$$

457 By the definition of \underline{t}_{k+1} , this satisfies

$$\begin{aligned} \|\tilde{\mathbf{u}} - \mathbf{u}_\alpha(t^*)\| &= \|\mathbf{u}_{\alpha, S_k^c}(t^*)\| \leq 4\rho + \|\mathbf{u}_{\alpha, S_k^c}(\bar{t}_k)\| \leq 4\rho + \underline{\epsilon}_k < 5\rho, \\ \|\tilde{\mathbf{v}} - \mathbf{v}_\alpha(t^*)\| &= \|\mathbf{v}_{\alpha, S_k^c}(t^*)\| \leq 4\rho + \|\mathbf{v}_{\alpha, S_k^c}(\bar{t}_k)\| \leq 4\rho + \underline{\epsilon}_k < 5\rho. \end{aligned}$$

458 Also

$$\|\tilde{\mathbf{u}} - \mathbf{u}^k\| + \|\tilde{\mathbf{v}} - \mathbf{s}^k \odot \mathbf{u}^k\| = \|\mathbf{u}_{\alpha, S_k}(t^*) - \mathbf{z}_{S_k}^k\| + \|\mathbf{v}_{\alpha, S_k}(t^*) - \mathbf{s}_{S_k}^k \odot \mathbf{z}_{S_k}^k\| \geq \rho' - 10\rho \geq \rho'/2.$$

459 Using (a) the strict minimum Assumption 4.3 with constant c_k , since $\|\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^k\| \leq \rho'$ and we take ρ'
460 small enough,

$$\begin{aligned} \mathcal{L}(\boldsymbol{\theta}_\alpha(t^*)) &\geq \mathcal{L}(\tilde{\boldsymbol{\theta}}) - 4\rho K_{R_k} \stackrel{(a)}{\geq} \mathcal{L}(\boldsymbol{\theta}^k) - 4\rho K_{R_k} + \frac{c_k(\rho')^2}{16} \\ &\geq \mathcal{L}(\boldsymbol{\theta}_\alpha(\bar{t}_k)) - (4\rho + \bar{\epsilon}_k)K_{R_k} + \frac{c_k(\rho')^2}{16} > \mathcal{L}(\boldsymbol{\theta}_\alpha(\bar{t}_k)). \end{aligned}$$

461 This is a contradiction because \mathcal{L} is nondecreasing along the gradient flow. \square

462 **Lemma A.6** (Log-scale approximation is correct during part (A)). *There are functions $\rho_{k+1}(\underline{\epsilon}_{k+1}) >$
463 0 and $\underline{\alpha}_{k+1}(\rho, \underline{\epsilon}_{k+1}) > 0$ such that for all $\rho < \rho_{k+1}$ and $\alpha < \underline{\alpha}_{k+1}$, and for all $t \in (\bar{t}_k, \underline{t}_{k+1})$ we
464 have for a constant C depending on k ,*

$$\|\log_\alpha(\mathbf{w}_\alpha(t)) - \mathbf{b}^k + (t - \bar{t}_k)\mathbf{g}(\boldsymbol{\theta}^k)\| < \rho\underline{\epsilon}_{k+1} + C\rho'(t - \bar{t}_k). \quad (26)$$

465 Furthermore, for all $i \in S_k^c$ and $t \in (\bar{t}_k, \underline{t}_{k+1})$ we have

$$\text{sgn}(g_i(\boldsymbol{\theta}_\alpha(t))) = \text{sgn}(g_i(\boldsymbol{\theta}^k)). \quad (27)$$

466 *Proof.* By Lemma A.5 and Lemma B.7, there is a constant C depending on $\boldsymbol{\theta}^k$ such that for all
467 $t \in (\bar{t}_k, \underline{t}_{k+1})$,

$$\|\mathbf{g}(\boldsymbol{\theta}_\alpha(t)) - \mathbf{g}(\boldsymbol{\theta}^k)\| \leq C\rho'.$$

468 For shorthand, write $\bar{\mathbf{g}}(\boldsymbol{\theta}^k) = \mathbf{g}(\boldsymbol{\theta}^k) + C\rho'\mathbf{1}$ and $\underline{\mathbf{g}}(\boldsymbol{\theta}^k) = \mathbf{g}(\boldsymbol{\theta}^k) - C\rho'\mathbf{1}$. Since $\mathbf{w}_\alpha(t) > 0$
469 entrywise as we have assumed without loss of generality (see Section A.3), we have the following
470 entrywise inequalities

$$\underline{\mathbf{g}}(\boldsymbol{\theta}^k) \odot \mathbf{w}_\alpha(t) < \mathbf{g}(\boldsymbol{\theta}_\alpha(t)) \odot \mathbf{w}_\alpha(t) < \bar{\mathbf{g}}(\boldsymbol{\theta}^k) \odot \mathbf{w}_\alpha(t). \quad (28)$$

471 Since the dynamics are given by $\frac{d\mathbf{w}_\alpha}{dt} = \log(1/\alpha)\mathbf{g}(\mathbf{w}_\alpha) \odot \mathbf{w}_\alpha$,

$$\mathbf{w}_\alpha(\bar{t}_k)e^{(t-\bar{t}_k)\log(1/\alpha)\underline{\mathbf{g}}(\boldsymbol{\theta}^k)} \leq \mathbf{w}_\alpha(t) \leq \mathbf{w}_\alpha(\bar{t}_k)e^{(t-\bar{t}_k)\log(1/\alpha)\bar{\mathbf{g}}(\boldsymbol{\theta}^k)}.$$

472 Taking the logarithms with base $\alpha \in (0, 1)$,

$$(t - \bar{t}_k)\underline{\mathbf{g}}(\mathbf{u}^k) \leq \log_\alpha(\mathbf{w}_\alpha(\bar{t}_k)) - \log_\alpha(\mathbf{w}_\alpha(t)) \leq (t - \bar{t}_k)\bar{\mathbf{g}}(\mathbf{u}^k).$$

473 The bound (26) follows since $\|\log_\alpha(\mathbf{w}_\alpha(\bar{t}_k)) - \mathbf{b}^k\| < \bar{\epsilon}_k < \rho\underline{\epsilon}_{k+1}$.

474 Finally, the claim (27) follows from (28) since $\text{sgn}(\bar{\mathbf{g}}(\boldsymbol{\theta}^k)) = \text{sgn}(\underline{\mathbf{g}}(\boldsymbol{\theta}^k)) = \text{sgn}(\mathbf{g}(\boldsymbol{\theta}^k))$ if we take
475 ρ small enough. \square

476 First, we show that the weights must move significantly by time roughly T_{k+1} . This is because of the
477 contribution of coordinate i_k .

478 **Lemma A.7** (\underline{t}_{k+1} is not much larger than T_{k+1}). *Suppose that $T_{k+1} < \infty$. Then there are
479 $\rho_{k+1}(\underline{\epsilon}_{k+1}) > 0$ and $\underline{\alpha}_{k+1}(\rho, \underline{\epsilon}_{k+1}) > 0$ such that for all $\rho < \rho_{k+1}$ and $\alpha < \underline{\alpha}_{k+1}$, the following
480 holds.*

$$\underline{t}_{k+1} < T_{k+1} + \underline{\epsilon}_{k+1}.$$

481 *Proof.* Assume by contradiction that $\underline{t}_{k+1} < T_{k+1} + \underline{\epsilon}_{k+1}$. For all times $t \in [\bar{t}_k, \min(\underline{t}_{k+1}, T_{k+1} +$
482 $\underline{\epsilon}_{k+1})]$, by Lemma A.6,

$$|\log_\alpha(w_{\alpha, i_k}(t)) - b_{i_k}^t + (t - \bar{t}_k)g_{i_k}(\boldsymbol{\theta}^k)| < O(\sqrt{\rho}).$$

483 Since we know $|\Delta_k(i_k) - (T_{k+1} - \bar{t}_k)| < \bar{\epsilon}_k$ and $b_{i_k}^t - \Delta_k(i_k)g_{i_k}(\boldsymbol{\theta}^k) \in \{0, 2\}$, it follows that

$$\log_\alpha(w_{\alpha, i_k}(T_{k+1} + \underline{\epsilon}_{k+1})) \notin (-|g_{i_k}(\boldsymbol{\theta}^k)|(\underline{\epsilon}_{k+1} - \bar{\epsilon}_{k+1}), 2 + |g_{i_k}(\boldsymbol{\theta}^k)|(\underline{\epsilon}_{k+1} - \bar{\epsilon}_{k+1})) + O(\sqrt{\rho}).$$

484 By taking ρ small enough, we see that $|g_{i_k}(\boldsymbol{\theta}^k)|(\underline{\epsilon}_{k+1} - \bar{\epsilon}_{k+1}) + O(\sqrt{\rho}) > \delta > 0$ for some $\delta > 0$
485 that is independent of α , so

$$\log_\alpha(w_{\alpha, i_k}(T_{k+1} + \underline{\epsilon}_{k+1})) \notin (-\delta, 2 + \delta).$$

486 So $|u_{\alpha, i_k}(T_{k+1} + \underline{\epsilon}_{k+1})| > 1$ by Lemma B.2. But by the construction of \underline{t}_{k+1} this means that
487 $\underline{t}_{k+1} < T_{k+1} + \underline{\epsilon}_{k+1}$. \square

488 Next, we show that until time \underline{t}_{k+1} , none of the coordinates in S_k^c move significantly, with the possible
489 exception of coordinate i_k .

490 **Lemma A.8** (No coordinates in $S_k^c \setminus \{i_k\}$ move significantly during part (A)). *Suppose $T_{k+1} < \infty$.
491 Then there are $\rho_{k+1}(\underline{\epsilon}_{k+1}) > 0$ and $\underline{\alpha}_{k+1}(\rho, \underline{\epsilon}_{k+1}) > 0$ such that for all $\rho < \rho_{k+1}$ and $\alpha < \underline{\alpha}_{k+1}$,
492 the following holds. There is a constant $c > 0$ depending on k such that for all $i \in S_k^c \setminus \{i_k\}$ and
493 $t \in [\bar{t}_k, \underline{t}_{k+1}]$,*

$$|u_{\alpha,i}(t) - u_{\alpha,i}(\bar{t}_k)|, |v_{\alpha,i}(t) - v_{\alpha,i}(\bar{t}_k)| < \alpha^c + \bar{\epsilon}_k.$$

Proof. The previous lemma combined with the inductive hypothesis gives

$$\underline{t}_{k+1} - \bar{t}_k < \Delta_k(i_k) + 2\underline{\epsilon}_{k+1} \setminus \{i_k\}.$$

494 We analyze the movement of each coordinate $i \in S_k^c \setminus \{i_k\}$ by breaking into two cases:

495 • Coordinate $i \neq i_k$ such that $b_i^k \in (0, 2)$. By Assumption 4.2, there is a unique winning
496 coordinate so $b_i^k - \tau g_i(\boldsymbol{\theta}^k) \in (c, 2 - c)$ for some constant $c > 0$ for all $\tau \in [0, \underline{t}_{k+1} - \bar{t}_k] \subseteq$
497 $[0, \Delta_k(i_k) + 2\underline{\epsilon}_{k+1}]$. By Lemma A.6, $\log_\alpha(u_{\alpha,i}(t)) \in (-c/2, 2 - c/2)$ for all times
498 $t \in [\bar{t}_k, \underline{t}_{k+1}]$. So by Lemma B.1, $|u_{\alpha,i}(t)|, |v_{\alpha,i}(t)| \leq \alpha^{c/4}$.

499 • Coordinate $i \neq i_k$ such that $b_i^k = 0$. By Lemma B.4, we must be in the corner case where
500 $i \in S_{k-1} \cap S_k^c$ (i.e., the coordinate was active in the previous stage but was dropped from
501 the support in this stage).

502 By Lemma B.4, since $b_i^k = 0$ we have $g_i(\boldsymbol{\theta}^k) < 0$. By Lemma A.6, this means
503 $\text{sgn}(g_i(\boldsymbol{\theta}_\alpha(t))) = \text{sgn}(g_i(\boldsymbol{\theta}^k)) < 0$ for all $t \in (\bar{t}_k, \underline{t}_{k+1})$.

504 We break the analysis into two parts. Since $b_i^k = 0$, the sign is $s_i^k = +1$. The inductive
505 hypothesis $\|\boldsymbol{\theta}_\alpha(\bar{t}_k) - \boldsymbol{\theta}^k\| < \bar{\epsilon}_k$ implies that $|u_{\alpha,i}(\bar{t}_k) - z_i^k| < \bar{\epsilon}_k$ and $|v_{\alpha,i}(\bar{t}_k) - z_i^k| < \bar{\epsilon}_k$.
506 For small enough $\bar{\epsilon}_k$ this means that $\text{sgn}(u_{\alpha,i}(\bar{t}_k)) = \text{sgn}(v_{\alpha,i}(\bar{t}_k)) = +1$. Now let
507 $t^* = \min(\underline{t}_{k+1}, \inf\{t > \bar{t}_k : v_{\alpha,i}(t) = 0\})$. Since $u_{\alpha,i}(t) > v_{\alpha,i}(t)$ without loss of
508 generality (see Section A.3), we have $\text{sgn}(u_{\alpha,i}(t)) = \text{sgn}(v_{\alpha,i}(t)) = +1$ for all $t \in [\bar{t}_k, t^*]$.
509 So $\frac{du_{\alpha,i}(t)}{dt}, \frac{dv_{\alpha,i}(t)}{dt} < 0$ for all $t \in [\bar{t}_k, t^*]$. So, for any $t \in [\bar{t}_k, t^*]$,

$$|u_{\alpha,i}(t) - u_{\alpha,i}(\bar{t}_k)|, |v_{\alpha,i}(t) - v_{\alpha,i}(\bar{t}_k)| < \bar{\epsilon}_k$$

510 Also, since $\log_\alpha(u_{\alpha,i}(t^*)) \approx 1$, by Lemma A.6 we have $t^* > c > 0$ for some constant c
511 independent of α . So for all $t \in [t^*, \underline{t}_{k+1}]$ we have $b_i^k - \tau g_i(\boldsymbol{\theta}^k) \in (c, 2 - c)$ for some
512 constant $c > 0$. So $|u_{\alpha,i}(t)|, |v_{\alpha,i}(t)| \leq \alpha^{c/4}$ for all $t \in [t^*, \underline{t}_{k+1}]$. The conclusion follows
513 by triangle inequality.

514 • Coordinate $i \neq i_k$ such that $b_i^k = 2$. The analysis is analogous to the case $b_i^k = 0$, except
515 that we have $s_i^k = -1$ instead and $g_i(\boldsymbol{\theta}^k) > 0$ by Lemma B.4.

516 □

517 Finally, we use this conclude that $\underline{t}_{k+1} \approx T_{k+1}$ and that the weights at coordinate i_k are the only
518 weights that change significantly, and by an amount approximately ρ .

519 **Lemma A.9** (Coordinate i_k wins the part (A) race at time $\underline{t}_{k+1} \approx T_{k+1}$). *Suppose that $T_{k+1} < \infty$.
520 Then there are $\rho_{k+1}(\underline{\epsilon}_{k+1}) > 0$ and $\underline{\alpha}_{k+1}(\rho, \underline{\epsilon}_{k+1}) > 0$ such that for all $\rho < \rho_{k+1}$ and $\alpha < \underline{\alpha}_{k+1}$,
521 the following holds.*

$$|\underline{t}_{k+1} - T_{k+1}| < \underline{\epsilon}_{k+1},$$

$$u_{\alpha,i_k}(\underline{t}_{k+1}) \in [\rho, 3\rho],$$

$$\text{sgn}(v_{\alpha,i_k}(\underline{t}_{k+1})) = s_{i_k}^{k+1}.$$

522 *Proof.* Let us analyze the case that $b_{i_k}^{k+1} \in (0, 2)$. Notice that $b_{i_k}^{k+1} = b_{i_k}^k - \Delta_k(i_k)g_{i_k}(\boldsymbol{\theta}^k) \in \{0, 2\}$
523 and that if $b_{i_k}^{k+1} = 0$ then $g_{i_k}(\boldsymbol{\theta}^k) > 0$ and if it is 2 then $b_{i_k}^{k+1} = g_{i_k}(\boldsymbol{\theta}^k) < 0$. So by Lemma A.6,
524 for all times $t \in [\bar{t}_k, \min(\underline{t}_{k+1}, T_{k+1} - \underline{\epsilon}_{k+1})]$, we have $w_{\alpha, i_k}(t) \in (c, 2 - c)$ for some $c > 0$. So for
525 small enough α by Lemma B.1, $|u_{\alpha, i_k}(t)|, |v_{\alpha, i_k}(t)| \leq \alpha^{c/2}$. Combining this with Lemma A.8, we
526 see that for $t \in [\bar{t}_k, \min(\underline{t}_{k+1}, T_{k+1} - \underline{\epsilon}_{k+1})]$ we have

$$\|\mathbf{u}_\alpha(t) - \mathbf{u}_\alpha(\bar{t}_k)\| + \|\mathbf{v}_\alpha(t) - \mathbf{v}_\alpha(\bar{t}_k)\| < 2(\alpha^c + \bar{\epsilon}_k)p < \rho,$$

527 for small enough α . So by definition of \underline{t}_{k+1} we must have $\underline{t}_{k+1} > T_{k+1} - \underline{\epsilon}_{k+1}$. Combined
528 with Lemma A.7, we conclude that $|T_{k+1} - \underline{t}_{k+1}| < \underline{\epsilon}_{k+1}$, which is the first claim of the lemma.
529 Furthermore, by Lemma A.8,

$$\sum_{i \in S_k^c \setminus \{i_k\}} |u_{\alpha, i}(\underline{t}_{k+1}) - u_{\alpha, i}(\bar{t}_k)| + |v_{\alpha, i}(\underline{t}_{k+1}) - v_{\alpha, i}(\bar{t}_k)| \leq 2p(\alpha^c + \bar{\epsilon}_k) < \rho/2,$$

530 so by definition of \underline{t}_{k+1} and triangle inequality we have $|u_{\alpha, i_k}(\underline{t}_{k+1})| + |v_{\alpha, i_k}(\underline{t}_{k+1})| \geq 4\rho - \rho/2 =$
531 $7\rho/2$. Also, since $u_{\alpha, i_k}^2(\underline{t}_{k+1}) - v_{\alpha, i_k}^2(\underline{t}_{k+1}) = \Theta(\alpha^2)$ we have $u_{\alpha, i_k}(\underline{t}_{k+1}) \in [\rho, 3\rho]$. Finally, if
532 $b_{i_k}^{k+1} = 2$, then $s_{i_k}^{k+1} = -1$ and $\log_\alpha(w_{\alpha, i_k}(\underline{t}_{k+1})) > 1.5$ so $\text{sgn}(v_{\alpha, i_k}(t)) < 0$ by Lemma B.3;
533 analogously, if $b_{i_k}^{k+1} = 0$, we have $s_{i_k}^{k+1} = 1$ and $\log_\alpha(w_{\alpha, i_k}(\underline{t}_{k+1})) < 0.5$ so $\text{sgn}(v_{\alpha, i_k}(\underline{t}_{k+1})) > 0$.

534 The case $b_{i_k}^k \in \{0, 2\}$ can be proved similarly to the analysis in Lemma A.8, where one shows that
535 during the first period of time the magnitudes of $|u_{i_k}(t)|$ and $|v_{i_k}(t)|$ decrease, until the sign of v_{i_k}
536 flips and they once again increase.

537 □

538 We have shown the claims (19), (20), (21) (22), and (23) for the time \underline{t}_{k+1} . In fact, if we let
539 $\underline{t}'_{k+1} \in [\bar{t}_k, \infty)$ be the first time t such that $u_{\alpha, i_k}(t) = \rho$ we still have (19), (20), (21) and (23) by the
540 same analysis as above, and (22) can be replaced with the slightly more convenient

$$u_{\alpha, i_k}(\underline{t}'_{k+1}) = \rho.$$

541 A.6.2 Analysis in case where $T_{k+1} = \infty$

542 In this case that T_{k+1} , we just have to show that the weights remain close to $\boldsymbol{\theta}^k$. We show that for
543 any $\underline{\epsilon}_{k+1} > 0$, there is $\underline{\alpha}_{k+1}(\underline{\epsilon}_{k+1}) > 0$ such that for all $\alpha < \underline{\alpha}_{k+1}$ and times $t \in [T_k + \underline{\epsilon}_{k+1}, T^*]$,

$$\|\boldsymbol{\theta}_\alpha(t) - \boldsymbol{\theta}^k\| < \underline{\epsilon}_{k+1}.$$

544 We can use Lemmas A.5 and A.6, which were developed for the case of $T_{k+1} < \infty$, but still hold for
545 $T_{k+1} = \infty$. Lemma A.5 guarantees that the weights do not move much until time \underline{t}_{k+1} , and so we
546 only need to show that $\underline{t}_{k+1} \geq T^*$ when we take ρ small enough. For this, observe that $g_i(\boldsymbol{\theta}^k) = 0$
547 for all $i \notin S_k$, because otherwise $T_{k+1} < \infty$. Therefore Lemma A.6 guarantees that until time
548 $\min(T^*, \underline{t}_{k+1})$ all weights are close to the original on the logarithmic scale. Namely,

$$\|\log_\alpha(\mathbf{w}_\alpha(t)) - \mathbf{b}^k\| < \rho\underline{\epsilon}_{k+1} + C\rho'(T^* - \bar{t}_k)$$

549 Furthermore, by the non-degeneracy Assumption 4.2 we know that $b_i^k \in (0, 2)$ for all $i \notin S_k$ by
550 Lemma B.4. So if we take ρ small enough and $\underline{\alpha}_{k+1}$ small enough, we must have that $\underline{t}_{k+1} \geq T^*$.

551 A.7 Dynamics from time \underline{t}_k to time \bar{t}_k (Nonlinear evolution for $O(1)$ unrescaled time)

552 Suppose that we know for some $k \leq K$ that for any $\underline{\epsilon}_k > 0$, there is $\rho_k(\underline{\epsilon}_k) > 0$ such that for all
553 $\rho < \rho_k$ there is $\underline{\alpha}_k(\rho, \underline{\epsilon}_k) > 0$ such that for all $\alpha < \underline{\alpha}_k$, there is a time $\underline{t}_k = \underline{t}_k(\alpha, \rho, \underline{\epsilon}_k)$ satisfying

$$|T_k - \underline{t}_k| < \underline{\epsilon}_k \tag{29}$$

$$\|\boldsymbol{\theta}_\alpha(\underline{t}_k) - \boldsymbol{\theta}^{k-1}\| < \underline{\epsilon}_k \tag{30}$$

$$\|\log_\alpha(\mathbf{w}_\alpha(\underline{t}_k)) - \mathbf{b}^k\| < \underline{\epsilon}_k \tag{31}$$

$$u_{\alpha, i_{k-1}}(\underline{t}_k) = \rho, \tag{32}$$

$$\text{sgn}(v_{\alpha, i_{k-1}}(\underline{t}_k)) = s_{i_{k-1}}^k. \tag{33}$$

554 Now we will show that for any $\bar{\epsilon}_k > 0$, there is $\bar{\alpha}_k = \bar{\alpha}_k(\bar{\epsilon}_k) > 0$ such that for all $0 < \alpha < \bar{\alpha}_k$, there
 555 is a time $\bar{t}_k = \bar{t}_k(\alpha, \bar{\epsilon}_k)$ satisfying

$$|T_k - \bar{t}_k| < \bar{\epsilon}_k \quad (34)$$

$$\|\boldsymbol{\theta}_\alpha(\bar{t}_k) - \boldsymbol{\theta}^k\| < \bar{\epsilon}_k \quad (35)$$

$$\|\log_\alpha(\mathbf{w}_\alpha(\bar{t}_k)) - \mathbf{b}^k\| < \bar{\epsilon}_k \quad (36)$$

556 We give the construction for \bar{t}_k . For any desired accuracy $\bar{\epsilon}_k > 0$ in this stage, we will construct an
 557 accuracy $\underline{\epsilon}_k = \underline{\epsilon}_k(\bar{\epsilon}_k) = \bar{\epsilon}_k/3 > 0$. We will also construct a $\rho = \rho(\underline{\epsilon}_k) > 0$ which is sufficiently small,
 558 and we will construct an cutoff for α equal to $\bar{\alpha}_k = \bar{\alpha}_{k+1}(\bar{\epsilon}_k) > 0$ which satisfies $\bar{\alpha}_k < \underline{\alpha}_k(\rho, \underline{\epsilon}_k)$.
 559 The values for these parameters $\underline{\epsilon}_k$ and ρ and $\bar{\alpha}_k$ will be chosen in the following lemma, and will
 560 depend only on $\bar{\epsilon}_k$.

561 **Lemma A.10** (New local minimum reached in time $O(1/\log(1/\alpha))$). *For any $\bar{\epsilon}_k > 0$, we can
 562 choose $\bar{\alpha}_k = \bar{\alpha}_k(\bar{\epsilon}_k) > 0$ small enough so that, for any $0 < \alpha < \bar{\alpha}_k$, there is $\bar{t}_k = \bar{t}_k(\alpha, \bar{\epsilon}_k)$ for
 563 which conditions (34) to (36) hold.*

564 *Furthermore, there is a constant C'' independent of α such that $|\boldsymbol{\theta}_\alpha(t)|/|\boldsymbol{\theta}_\alpha(\bar{t}_k)| \in [1/C'', C'']^{2p}$ at
 565 all times $t \in [\bar{t}_k, \bar{t}_k]$.*

566 *Proof.* Let $\underline{t}_k = \underline{t}_k(\alpha, \rho, \underline{\epsilon}_k)$ be given by the induction. Let us compare the dynamics starting at
 567 $\boldsymbol{\theta}_\alpha(\underline{t}_k)$ with the dynamics starting at $\tilde{\boldsymbol{\theta}}(\underline{t}_k) = (\tilde{u}(\underline{t}_k), \tilde{v}(\underline{t}_k))$ which is given by

$$\tilde{u}_i(\underline{t}_k) = \begin{cases} u_{\alpha, i}(\underline{t}_k), & i \in S_{k-1} \cup \{i_{k-1}\} \\ 0, & \text{otherwise} \end{cases} \quad \text{and} \quad \tilde{v}_i(\underline{t}_k) = \begin{cases} v_{\alpha, i}(\underline{t}_k), & i \in S_{k-1} \cup \{i_{k-1}\} \\ 0, & \text{otherwise} \end{cases}$$

568 and run with

$$\frac{d\tilde{\boldsymbol{\theta}}}{dt} = -\log(1/\alpha)\nabla_{\mathbf{w}}\mathcal{L}(\tilde{\boldsymbol{\theta}}).$$

569 By Assumption 4.4 we know there exists a unique solution $\tilde{\boldsymbol{\theta}} : [\underline{t}_k, \infty) \rightarrow \mathbb{R}^p$ as long as we take $\underline{\epsilon}_k$
 570 small enough because $\text{supp}(\tilde{\boldsymbol{\theta}}(\underline{t}_k)) = S_{k-1} \cup \{i_{k-1}\}$ and $\|\tilde{\boldsymbol{\theta}}_i(\underline{t}_k) - \boldsymbol{\theta}^{k-1}\| < \underline{\epsilon}_k$. Furthermore, by
 571 Assumption 4.4 if we take $\underline{\epsilon}_k$ small enough there must be a time $\tau := \tau(\bar{\epsilon}_k, \rho) < \infty$ such that

$$\|\tilde{\boldsymbol{\theta}}(t) - \boldsymbol{\theta}^k\| < \bar{\epsilon}_k/2 \text{ for } t \geq \underline{t}_k + \tau/\log(1/\alpha) \quad (37)$$

572 Define

$$\bar{t}_k = \underline{t}_k + \tau/\log(1/\alpha).$$

573 So for α small enough, $|T_k - \bar{t}_k| < 2\underline{\epsilon}_k < \bar{\epsilon}_k$, proving (34).

574 We now compare $\boldsymbol{\theta}_\alpha(\bar{t}_k)$ with $\tilde{\boldsymbol{\theta}}(\bar{t}_k)$, and show that if we take α small enough, then the dynamics of $\tilde{\boldsymbol{\theta}}$
 575 closely match the dynamics of $\boldsymbol{\theta}_\alpha(t)$ for times $\underline{t}_k + O(1/\log(1/\alpha))$. The argument uses Gronwall's
 576 inequality. Let $t^* = \inf\{t > \underline{t}_k : \|\tilde{\boldsymbol{\theta}}(t^*) - \boldsymbol{\theta}_\alpha(t^*)\| > 1/3\}$. For times $t \in [\underline{t}_k, t^*)$ by Lemma B.7 we
 577 have

$$\left\| \frac{d}{dt}\tilde{\boldsymbol{\theta}}(t) - \frac{d}{dt}\boldsymbol{\theta}_\alpha(t) \right\| = \log(1/\alpha)\|\nabla_{\boldsymbol{\theta}}\mathcal{L}(\tilde{\boldsymbol{\theta}}(t)) - \nabla_{\boldsymbol{\theta}}\mathcal{L}(\boldsymbol{\theta}_\alpha(t))\| \leq K_{\tilde{\boldsymbol{\theta}}(t)}\log(1/\alpha)\|\tilde{\boldsymbol{\theta}}(t) - \boldsymbol{\theta}_\alpha(t)\|,$$

578 where $K_{\tilde{\boldsymbol{\theta}}(t)}$ is the smoothness constant from Lemma B.7. Note that since $\|\tilde{\boldsymbol{\theta}}(t)\| < \infty$ for large
 579 enough t by (37), the trajectory of $\tilde{\boldsymbol{\theta}}$ must lie in a compact set. Therefore, there must be a finite
 580 set of times $s_1, \dots, s_m \in [\underline{t}_k, t^*)$ such that $\cup_{t \in [\underline{t}_k, t^*)} B(\tilde{\boldsymbol{\theta}}(t), 1/2) \subseteq \cup_{i=1}^m B(\tilde{\boldsymbol{\theta}}(s_i), 3/4)$. So letting
 581 $C = \max_{i=1}^m K_{\tilde{\boldsymbol{\theta}}(s_i)} < \infty$ for all times $t \in [\underline{t}_k, t^*)$ we have

$$\frac{d}{dt}\|\tilde{\boldsymbol{\theta}}(t) - \boldsymbol{\theta}_\alpha(t)\| \leq C\log(1/\alpha)\|\tilde{\boldsymbol{\theta}}(t) - \boldsymbol{\theta}_\alpha(t)\|.$$

582 By Gronwall's inequality, for all times $t \in [\underline{t}_k, t^*)$,

$$\|\tilde{\boldsymbol{\theta}}(t) - \boldsymbol{\theta}_\alpha(t)\| \leq \|\tilde{\boldsymbol{\theta}}(\underline{t}_k) - \boldsymbol{\theta}_\alpha(\underline{t}_k)\| \exp(C\log(1/\alpha)(t - \underline{t}_k)).$$

583 We know from Lemma A.8 that there is a constant $c > 0$ such that for any small enough $0 < \alpha < \underline{\alpha}_k$,
584 such that

$$\|\tilde{\boldsymbol{\theta}}(\underline{t}_k) - \boldsymbol{\theta}_\alpha(\underline{t}_k)\| < \alpha^c$$

585 If we take α small enough that $\alpha^c \exp(C\tau) < \bar{\epsilon}_k/2 < 1/3$, we must have $t^* > \underline{t}_k + \tau/\log(1/\alpha)$
586 and so we prove (35)

$$\|\boldsymbol{\theta}^k - \boldsymbol{\theta}_\alpha(\bar{t}_k)\| \leq \bar{\epsilon}_k/2 + \|\tilde{\boldsymbol{\theta}}(\bar{t}_k) - \boldsymbol{\theta}_\alpha(\bar{t}_k)\| < \bar{\epsilon}_k.$$

587 It remains to show that (36) is satisfied. Since $\|\tilde{\boldsymbol{\theta}}(t) - \boldsymbol{\theta}_\alpha(t)\| < 1/3$ for all $t \in [\underline{t}_k, \bar{t}_k]$, it holds that
588 the trajectory of $\boldsymbol{\theta}_\alpha(t)$ lies in a compact set. So by Lemma B.7 we have $\|\mathbf{g}(\boldsymbol{\theta}_\alpha(t))\| < C'$ for some
589 constant C' at all times $t \in [\underline{t}_k, \bar{t}_k]$. Since $\frac{1}{\log(1/\alpha)} \left| \frac{dw_{\alpha,i}}{dt} \right| = |w_{\alpha,i}(t)| |g_i(\mathbf{w}_\alpha(t))| < C' |w_{\alpha,i}(t)|$,
590 we must have $|w_{\alpha,i}(t)|/|w_{\alpha,i}(\underline{t}_k)| \in [1/C'', C'']$ for some constant C'' independent of α and all
591 $t \in [\underline{t}_k, \bar{t}_k]$. Therefore, (36) follows from (31). A similar argument shows that $|\boldsymbol{\theta}_\alpha(t)/\boldsymbol{\theta}_\alpha(\underline{t}_k)| \in$
592 $[1/C''', C''']^{2p}$.

593 □

594 A.8 Concluding the proof of Theorem A.4

595 We have shown that Theorem 4.5 is true for solutions $\boldsymbol{\theta}_\alpha : [0, T^*] \rightarrow \mathbb{R}^{2p}$ to the gradient flow,
596 where $T_* \in (T_K, T_{K+1})$. To establish Theorem A.4 it remains only to show that for any $T_* \in$
597 (T_K, T_{K+1}) and small enough α such a solution to the gradient flow exists and is unique. To
598 see this, note that in the inductive proof of the invariants we construct a sequence of times $0 =$
599 $\bar{t}_0 \leq \underline{t}_1 \leq \bar{t}_1 \leq \dots \leq \bar{t}_K \leq \underline{t}_{K+1} > T_*$, where we guarantee that any gradient flow solution
600 $\boldsymbol{\theta}_\alpha : [0, \underline{t}_{k+1}] \rightarrow \mathbb{R}^p$ satisfies $\boldsymbol{\theta}_\alpha \in \cup_{k \in \{0, \dots, K\}} B(\boldsymbol{\theta}^k, 1)$ for all $t \in \cup_{k \in \{0, \dots, K\}} [\bar{t}_k, \underline{t}_{k+1}]$. And also
601 for $t \in \cup_{k \in \{0, \dots, K-1\}} [\underline{t}_k, \bar{t}_{k+1}]$, we have $\boldsymbol{\theta}_\alpha(t) \in B(0, C''_k \boldsymbol{\theta}^k)$ for some constant C''_k independent
602 of α by Lemma A.10. So $\boldsymbol{\theta}_\alpha(t) \in B(0, C_K)$ for some constant C_K at all times $t \in [0, T^*]$. By
603 Lemma B.7, the loss gradient $\nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}) = (\mathbf{v} \odot \mathbf{g}(\boldsymbol{\theta}), \mathbf{u} \odot \mathbf{g}(\boldsymbol{\theta}))$ is Lipschitz-continuous on the
604 compact set $B(0, C_K)$. So $\boldsymbol{\theta}_\alpha : [0, T^*] \rightarrow \mathbb{R}^p$ exists and is unique by the Cauchy-Lipschitz theorem.

605 □

606 B Technical lemmas

607 B.1 Relating the sum of the weights to the original weights using the conservation law

608 **Lemma B.1.** *If for some constant $0 < c < 1$ we have $\log_\alpha(w_{\alpha,i}(t)) \in (c, 2 - c)$, then for small*
609 *enough α*

$$\max(|u_{\alpha,i}(t)|, |v_{\alpha,i}(t)|) \leq \alpha^{c/2}.$$

610 *Proof.* Let $\tilde{\mathbf{w}}_\alpha(t) = \mathbf{u}_\alpha(t) - \mathbf{v}_\alpha(t)$. By the conservation law (5), $w_{\alpha,i}(t)\tilde{w}_{\alpha,i}(t) =$
611 $w_{\alpha,i}(0)\tilde{w}_{\alpha,i}(0) = u_{\alpha,i}(0)^2 - v_{\alpha,i}(0)^2$. By the non-degeneracy of initialization (Assumption 4.2),
612 the right-hand-side is $\Theta(\alpha^2)$. So if $\log_\alpha(w_{\alpha,i}(t)) \in (c, 2 - c)$ then for small enough α , we
613 have $\log_\alpha(|\tilde{w}_{\alpha,i}(t)|) \in (3c/4, 2 - 3c/4)$. So $|u_{\alpha,i}(t)| \leq |w_{\alpha,i}(t) + \tilde{w}_{\alpha,i}(t)| \leq \alpha^{c/2}$ and
614 $|v_{\alpha,i}(t)| \leq |w_{\alpha,i}(t) - \tilde{w}_{\alpha,i}(t)| \leq \alpha^{c/2}$. □

615 **Lemma B.2.** *If for some constant $0 < c$ we have $\log_\alpha(w_{\alpha,i}(t)) \notin (-c, 2 + c)$, then for small enough*
616 *α ,*

$$|u_{\alpha,i}(t)| > 1.$$

617 *Proof.* Define $\tilde{\mathbf{w}}_\alpha = \mathbf{u}_\alpha - \mathbf{v}_\alpha$ as in the proof of Lemma B.1. If $\log_\alpha(w_{\alpha,i}(t)) < -c$ then
618 $\log_\alpha(|\tilde{w}_{\alpha,i}(t)|) > 2 - c/2$ for small enough α , so $|u_i(t)| > \alpha^{-c} - \alpha^{2-c/2} > 1$. Similarly, if
619 $\log_\alpha(w_{\alpha,i}(t)) > 2 + c$ then $\log_\alpha(|\tilde{w}_{\alpha,i}(t)|) < -c/2$ so $|u_i(t)| > \alpha^{-c/2} - \alpha^{2+c} > 1$. □

620 **Lemma B.3.** *If for some constant $c > 0$, there is small enough α such that if we have $\log_\alpha(w_{\alpha,i}(t)) >$
621 $1 + c$ then $\text{sgn}(v_{\alpha,i}(t)) < 0$. Otherwise, if $\log_\alpha(w_{\alpha,i}(t)) < 1 - c$ then $\text{sgn}(v_{\alpha,i}(t)) > 0$.*

622 *Proof.* Follows from $\mathbf{v}_\alpha = \frac{1}{2}(\mathbf{w}_\alpha - \tilde{\mathbf{w}}_\alpha)$. Recall that $w_\alpha(t) > 0$ and notice that $\tilde{w}_\alpha(t) > 0$.
623 In the first case, $w_{\alpha,i}(t) < \alpha^{1+c}$ and $\tilde{w}_{\alpha,i}(t) > \alpha^{1-c/2}$. In the latter case $w_{\alpha,i}(t) > \alpha^{1-c}$ and
624 $\tilde{w}_{\alpha,i}(t) < \alpha^{1+c/2}$. \square

625 B.2 Sign of gradients on coordinates that leave support

626 **Lemma B.4.** *For any $k \geq 1$ and $i \in S_k^c$, if $b_i^k \in \{0, 2\}$ then we must have $i \in \text{supp}(\mathbf{u}^{k-1}) \setminus$
627 $\text{supp}(\mathbf{u}^k)$, and we must have $g_i(\mathbf{u}^k) < 0$ if $b_i^k = 0$ and $g_i(\boldsymbol{\theta}^k) > 0$ if $b_i^k = 2$. In particular,
628 $\Delta_k(i_k) > 0$ for all k .*

629 *Proof.* This is by induction on k and using the non-degeneracy Assumption 4.2. \square

630 B.3 Local Lipschitzness and smoothness

631 We provide several technical lemmas on the local Lipschitzness and smoothness of ℓ , h , and \mathbf{g} .

632 **Lemma B.5.** *The function $\ell(\mathbf{y}, \cdot)$ is locally Lipschitz and smooth in its second argument: for any
633 $R > 0$, there exists K_R such that for any $\boldsymbol{\zeta}, \boldsymbol{\zeta}' \in B(0, R)$*

$$\begin{aligned} |\ell(\mathbf{y}, \boldsymbol{\zeta}) - \ell(\mathbf{y}, \boldsymbol{\zeta}')| &\leq K_R \|\boldsymbol{\zeta} - \boldsymbol{\zeta}'\| \\ \|D\ell(\mathbf{y}, \boldsymbol{\zeta}) - D\ell(\mathbf{y}, \boldsymbol{\zeta}')\| &\leq K_R \|\boldsymbol{\zeta} - \boldsymbol{\zeta}'\|, \end{aligned}$$

634 *almost surely over \mathbf{y} . Here $D\ell(\mathbf{y}, \cdot)^\top \in \mathbb{R}^{d_{out}}$ is the derivative in the second argument.*

635 *Proof.* Since ℓ is continuously twice-differentiable, for each $\mathbf{y} \in \mathbb{R}^{d_y}$, $\boldsymbol{\zeta} \in \mathbb{R}^{d_{out}}$ there is $K_{\mathbf{y}, \boldsymbol{\zeta}} < \infty$
636 such that for all $\mathbf{y} \in B(\mathbf{y}, 1/K_{\mathbf{y}, \boldsymbol{\zeta}})$ and $\boldsymbol{\zeta}' \in B(\boldsymbol{\zeta}, 1/K_{\mathbf{y}, \boldsymbol{\zeta}})$ we have

$$\|D\ell(\mathbf{y}', \boldsymbol{\zeta}')\| \leq K_{\mathbf{y}, \boldsymbol{\zeta}} \quad \text{and} \quad \|D^2\ell(\mathbf{y}', \boldsymbol{\zeta}')\| \leq K_{\mathbf{y}, \boldsymbol{\zeta}},$$

637 where $D\ell$ and $D^2\ell$ denote the first and second derivative in the second argument. So for all such
638 $\mathbf{y}' \in B(\mathbf{y}, 1/K_{\mathbf{y}, \boldsymbol{\zeta}})$ and $\boldsymbol{\zeta}', \boldsymbol{\zeta}'' \in B(\boldsymbol{\zeta}, 1/K_{\mathbf{y}, \boldsymbol{\zeta}})$ we have

$$|\ell(\mathbf{y}', \boldsymbol{\zeta}') - \ell(\mathbf{y}', \boldsymbol{\zeta}'')| \leq K_{\mathbf{y}, \boldsymbol{\zeta}} \|\boldsymbol{\zeta}' - \boldsymbol{\zeta}''\| \quad \text{and} \quad |D\ell(\mathbf{y}', \boldsymbol{\zeta}') - D\ell(\mathbf{y}', \boldsymbol{\zeta}'')| \leq K_{\mathbf{y}, \boldsymbol{\zeta}} \|\boldsymbol{\zeta}' - \boldsymbol{\zeta}''\|.$$

639 Cover the set $\{(\mathbf{y}, \boldsymbol{\zeta}) : \|\mathbf{y}\| \leq C, \|\boldsymbol{\zeta}\| \leq R\}$ with the balls $\cup_{\mathbf{y}} B(\mathbf{y}, 1/K_{\mathbf{y}, \boldsymbol{\zeta}})$. By compactness,
640 there is a finite subcover $(\mathbf{y}_1, \boldsymbol{\zeta}_1), \dots, (\mathbf{y}_r, \boldsymbol{\zeta}_r)$, so we can take $K_R = \max_{i \in [r]} K_{\mathbf{y}_i, \boldsymbol{\zeta}_i} < \infty$ and the
641 lemma holds since $\|\mathbf{y}\| \leq C$ almost surely by Assumption 2.1. \square

642 **Lemma B.6.** *The function $h(\mathbf{x}; \cdot)$ is locally bounded, Lipschitz and smooth in its second argument:
643 for any $R > 0$ there exists K_R such that for any $\boldsymbol{\psi}, \boldsymbol{\psi}' \in B(0, R)$,*

$$\begin{aligned} \|h(\mathbf{x}; \boldsymbol{\psi})\| &\leq K_R \\ \|h(\mathbf{x}; \boldsymbol{\psi}) - h(\mathbf{x}; \boldsymbol{\psi}')\| &\leq K_R \|\boldsymbol{\psi} - \boldsymbol{\psi}'\| \\ \|Dh(\mathbf{x}; \boldsymbol{\psi}) - Dh(\mathbf{x}; \boldsymbol{\psi}')\| &\leq K_R \|\boldsymbol{\psi} - \boldsymbol{\psi}'\|, \end{aligned}$$

644 *almost surely over \mathbf{x} . Here $Dh(\mathbf{x}, \cdot) \in \mathbb{R}^{d_{out}} \times \mathbb{R}^p$ is the derivative in the second argument.*

645 *Proof.* Analogous to proof of Lemma B.5, using continuous twice-differentiability of h and bounded-
646 ness of $\|\mathbf{x}\|$. \square

647 **Lemma B.7** (Local Lipschitzness of loss and loss derivative). *When $\boldsymbol{\theta} = (\mathbf{u}, \mathbf{v}) \in \mathbb{R}^{2p}$ and
648 $f_{\text{NN}}(\mathbf{x}; \boldsymbol{\theta}) = h(\mathbf{x}; \mathbf{u} \odot \mathbf{v})$ the following holds for $\mathbf{g}(\boldsymbol{\theta})$ defined in (4). For any $R > 0$, there exists
649 $K_R < \infty$ such that for any $\boldsymbol{\theta}, \boldsymbol{\theta}' \in B(0, K_R)$,*

$$\begin{aligned} \|\mathbf{g}(\boldsymbol{\theta}) - \mathbf{g}(\boldsymbol{\theta}')\| &\leq K_R \|\boldsymbol{\theta} - \boldsymbol{\theta}'\| \\ \|\nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}) - \nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}')\| &\leq K_{\boldsymbol{\theta}} \|\boldsymbol{\theta} - \boldsymbol{\theta}'\| \\ |\mathcal{L}(\boldsymbol{\theta}) - \mathcal{L}(\boldsymbol{\theta}')| &\leq K_R \|\boldsymbol{\theta} - \boldsymbol{\theta}'\|. \end{aligned}$$

650 *Proof.* Let $\boldsymbol{\theta} = (\mathbf{u}, \mathbf{v}), \boldsymbol{\theta}' = (\mathbf{u}', \mathbf{v}')$. This follows immediately from the local Lipschitzness and
651 smoothness of h and ℓ in Lemmas B.5 and B.6, as well as

$$\|\mathbf{g}(\boldsymbol{\theta}) - \mathbf{g}(\boldsymbol{\theta}')\| = \|\mathbb{E}_{\mathbf{x}, \mathbf{y}} [Dh(\mathbf{x}; \mathbf{u} \odot \mathbf{v})^\top D\ell(\mathbf{y}, h(\mathbf{x}; \mathbf{u} \odot \mathbf{v}))^\top - Dh(\mathbf{x}; \mathbf{u}' \odot \mathbf{v}')^\top D\ell(\mathbf{y}, h(\mathbf{x}; \mathbf{u}' \odot \mathbf{v}'))^\top]\|.$$

652 \square

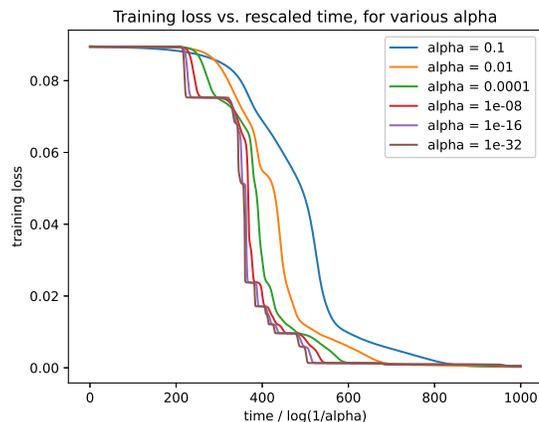


Figure 6: Evolution of loss versus rescaled time initializing at various scalings α in the toy task of learning an attention head with diagonal weights. The loss curves converge as $\alpha \rightarrow 0$ to a curve with loss plateaus and sharp decreases, as predicted by the theory.

653 C Experimental validation of the assumptions in Theorem 4.5

654 In Figures 6, 7, and 8, we plot the evolution of the losses, of the entries of $\mathbf{W}_K \mathbf{W}_Q^\top =$
 655 $\text{diag}(\mathbf{w}_K) \text{diag}(\mathbf{w}_Q)$, and of the entries of $\mathbf{W}_V \mathbf{W}_O^\top = \text{diag}(\mathbf{w}_V) \text{diag}(\mathbf{w}_O)$ in the toy task of
 656 training an attention head (1) with diagonal weights. The model is trained with SGD on the mean-
 657 squared error loss on 1000 random samples (\mathbf{X}, \mathbf{y}) . Each random sample has $\mathbf{X} \in \mathbb{R}^{10 \times 50}$, which a
 658 sequence of 10 tokens, each of dimension 50, which is distributed as isotropic Gaussians. The label
 659 \mathbf{y} is given by a randomly-generated teacher model that is also an attention head (1) with diagonal
 660 weights. In Figures 6, 7, and 8, for $\alpha \in \{0.1, 0.01, 0.0001, 10^{-8}, 10^{-16}, 10^{-32}\}$ we plot the evolu-
 661 tion of the loss and of the weights when initialized at $\boldsymbol{\theta}(0) = \alpha \boldsymbol{\theta}_0$, for some random Gaussian $\boldsymbol{\theta}_0$.
 662 Qualitatively, as $\alpha \rightarrow 0$ we observe that the loss curve and the trajectories of the weights appear to
 663 converge to a limiting stagewise dynamics, where there are plateaus followed by movement on short
 664 time-scales, as predicted by the theory.

665 **Validation of Assumption 4.2 (non-degeneracy of dynamics)** As $\alpha \rightarrow 0$, notice that the stages
 666 appear to separate and happen at distinct times. Furthermore, at no stage do any of the nonnegligible
 667 coordinates leave the support of $\boldsymbol{\theta}$, so the extra technical condition on coordinates $i \in \text{supp}(\boldsymbol{\theta}^k) \setminus$
 668 $\text{supp}(\boldsymbol{\theta}^{k-1})$ in Assumption 4.2 is automatically satisfied since $\text{supp}(\boldsymbol{\theta}^k) \setminus \text{supp}(\boldsymbol{\theta}^{k-1})$ is empty.

669 **Validation of Assumption 4.3 (stationary points are strict local minima)** In Figure 9 we consider
 670 the $\alpha = 10^{-32}$ trajectory, since this is closest to the dynamics in the $\alpha \rightarrow 0$ limit. We randomly select
 671 several epochs. Since the transitions between stages are a vanishing fraction of the total training time,
 672 each of these randomly-selected epochs is likely during a plateau, as we see in the figure. For each
 673 epoch perform the following experiment. For each nonnegligible coordinate of the weights (those
 674 where the weight is of magnitude greater than the threshold $\tau = 10^{-5}$), we perturb the weights by
 675 adding noise of standard deviation 0.05. We then run the training dynamics starting at this perturbed
 676 initialization for 1000 epochs. We observe that the training dynamics quickly converge to the original
 677 unperturbed initialization, indicating that the weights were close to a strict local minimum of the loss.

678 **Validation of Assumption 4.4 (noise-robustness of dynamics)** In Figure 10 we perform the same
 679 experiment as in Figure 9, except that the epochs we select to perturb the weights are those where
 680 there is a newly-nonnegligible coordinate (less than 10^{-5} in magnitude in the previous epoch, and
 681 more than 10^{-5} in magnitude in this epoch). We find that the nonlinear dynamics are robust and tend
 682 to the limiting endpoint even under a random Gaussian perturbation of standard deviation 10^{-2} on
 683 each of the nonnegligible coordinates, supporting Assumption 4.4.

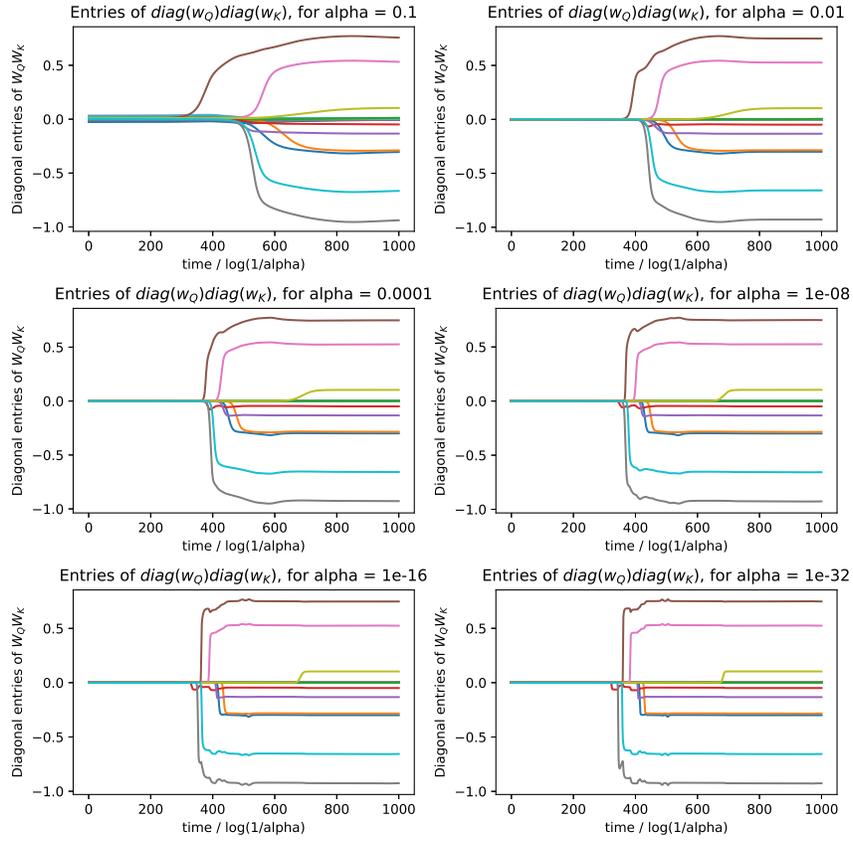


Figure 7: Evolution of $\text{diag}(\mathbf{w}_Q)\text{diag}(\mathbf{w}_K)$ entries over rescaled time initializing at various scalings α . Notice that as $\alpha \rightarrow 0$, the training trajectories tend to a limiting trajectory. Each line corresponds to a diagonal entry of $\text{diag}(\mathbf{w}_Q)\text{diag}(\mathbf{w}_K)$.

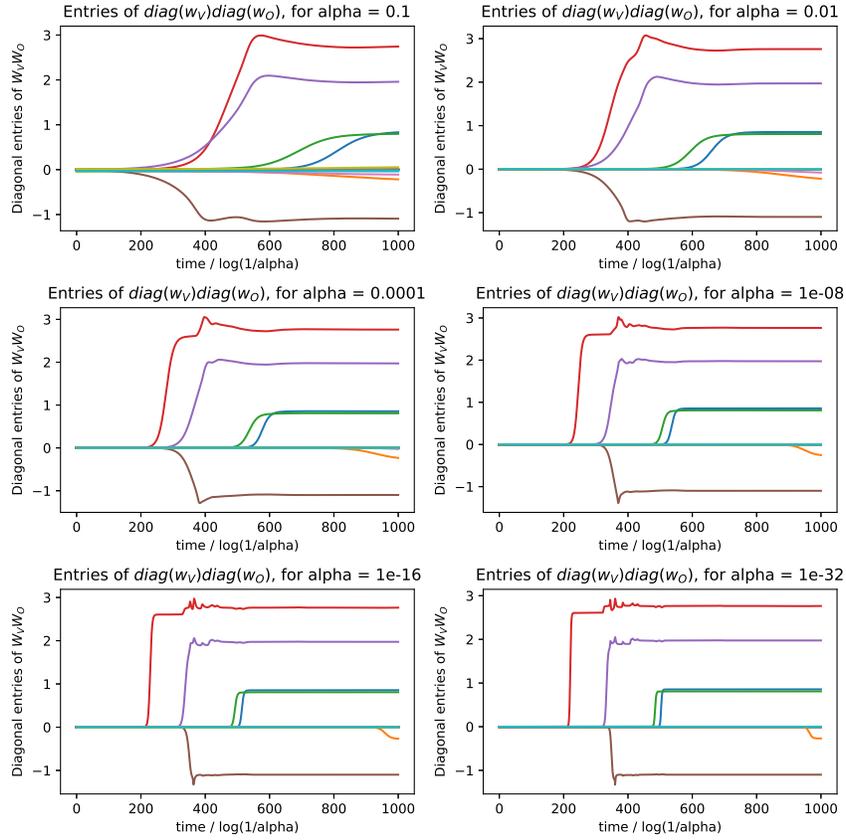


Figure 8: Evolution of $\text{diag}(\mathbf{w}_V)\text{diag}(\mathbf{w}_O)$ entries in the toy task of learning an attention head with diagonal weights. Each line corresponds to the evolution of an entry of $\text{diag}(\mathbf{w}_V)\text{diag}(\mathbf{w}_O)$ over rescaled time. Each plot corresponds to a different initialization magnitude α . Notice that as $\alpha \rightarrow 0$, the training trajectories tend to a limiting trajectory.

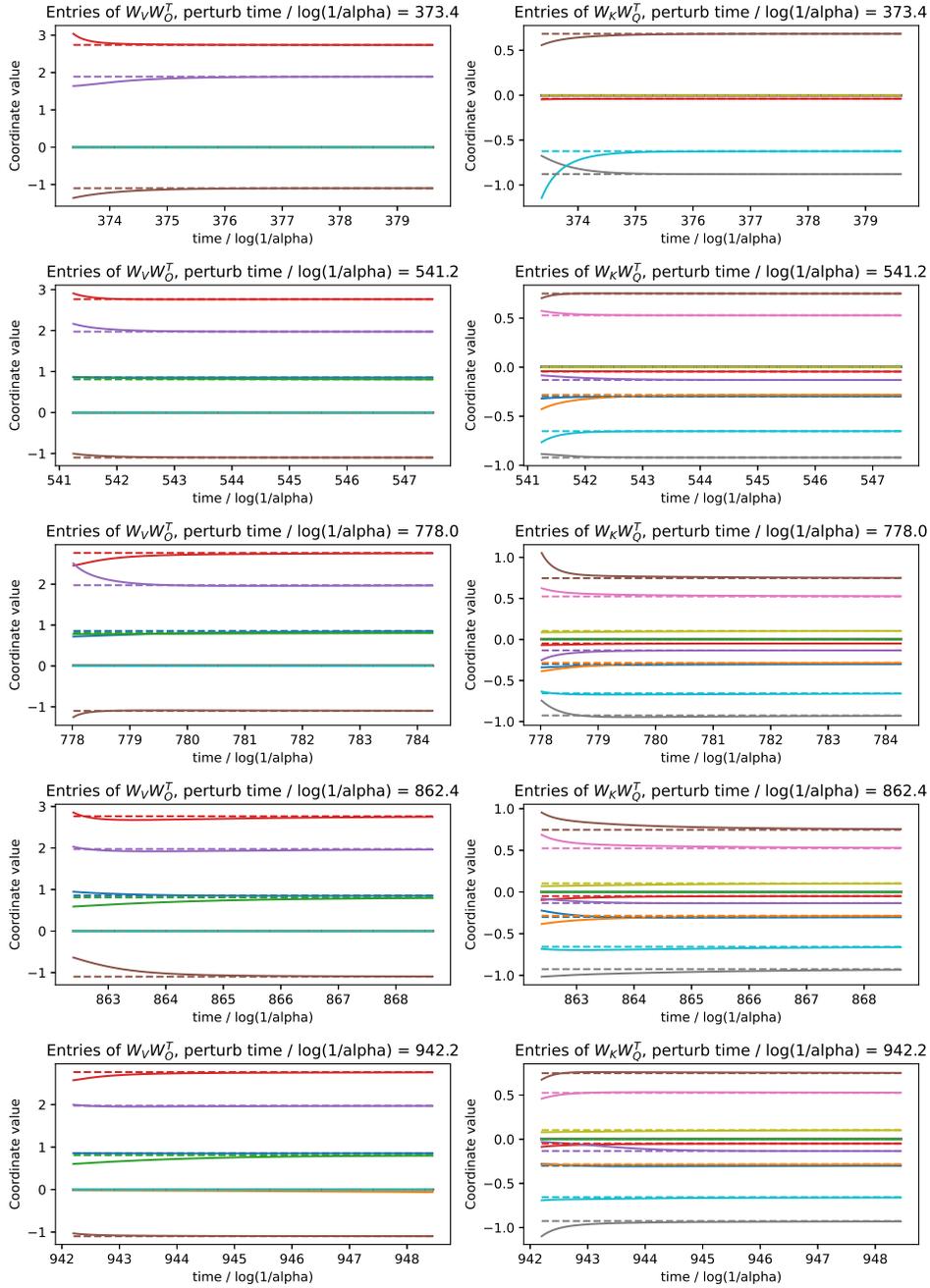


Figure 9: Evolution of weights of toy attention model under perturbation, validating Assumption 4.3. At 5 different random times during training, we perturb the nonnegligible weight coordinates and continue to train with SGD. The evolution of each of the weights under the initial perturbation (solid line) is compared to the original evolution without perturbation (dashed line). Observe that the training dynamics quickly brings each weight back to the unperturbed weight trajectory, indicating that the weights are originally close to a strict local minimum.

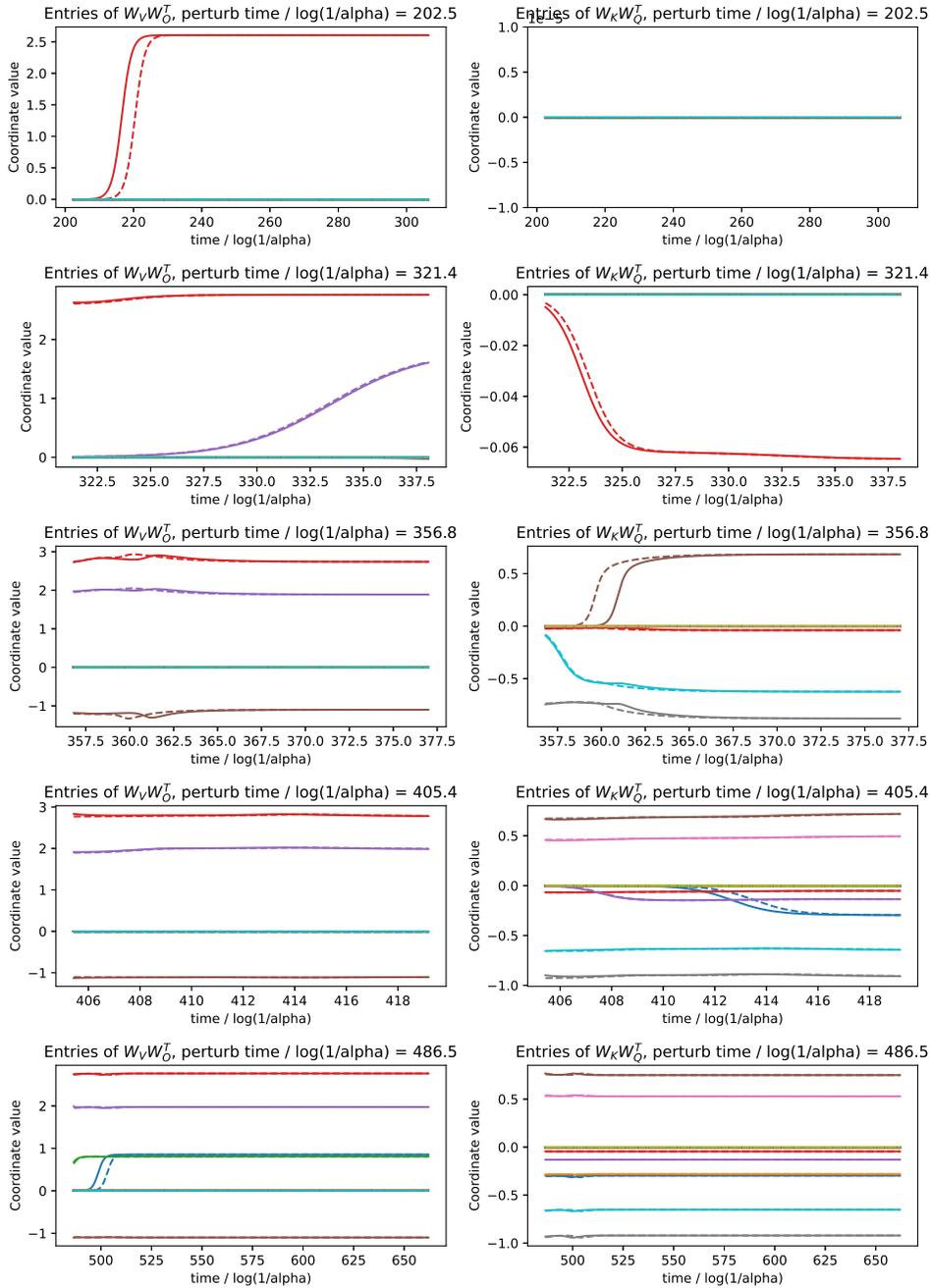


Figure 10: Validating Assumption 4.4 with the same experiment as in Figure 9, except that the epochs for the perturbation chosen are those where there is a newly nonnegligible coordinate. Perturbed dynamics (solid lines) are again robust to perturbation and track the original dynamics (dashed lines).

684 **D Vision Transformers**

685 The practice of training transformer models often deviate substantially from the assumptions made
686 in our theoretical analysis, and it is unclear to what extent gradual rank increase behaviour, and
687 a low rank bias are manifested in setups more common in practical applications. To gauge the
688 relevancy of our findings we conduct experiments on popular vision benchmarks, using algorithms
689 and hyperparameters common in the literature. We use the stable rank given by $\frac{\|s\|_F^2}{\|s\|_2^2}$, where s is
690 the spectrum, as a smooth approximation of rank. We track the value of the stable rank for the
691 different attention matrices throughout training. Although we do not expect our theoretical results to
692 hold precisely in practice, we find evidence of gradual increase in stable rank, leading to a low
693 rank bias Figures 12, 14 and 16. In these experiments we use off the shelf vision transformers (ViT)
694 Dosovitskiy et al. (2020) trained on popular vision benchmarks. For the Cifar-10/100 datasets we
695 use a ViT with 6 layers, patchsize of 4, 8 heads per self attention layer, an embedding and MLP
696 dimension of 512, and a head dimension of 128. We train the model using the Adam optimizer for 500
697 epochs with a base learning rate of 1e-4, a cyclic learning rate decay with a linear warmup schedule
698 for 15 epochs and a batchsize of 512. For Imagenet, we use the ViT-Base/16 from Dosovitskiy et al.
699 (2020) trained with Adam for 360 epochs with a base learning rate of 3e-3, a cyclic learning rate
700 decay with a linear warmup schedule for 15 epochs and a batchsize of 4096. We use no weight
701 decay or dropout in our experiments. All models were initialized using the default initialization scale.
702 Our results are summarized in Figures 11 and 12 for Cifar-10, Figures 13 and 14 for Cifar-100 and
703 Figures 15 and 16 for imagenet.

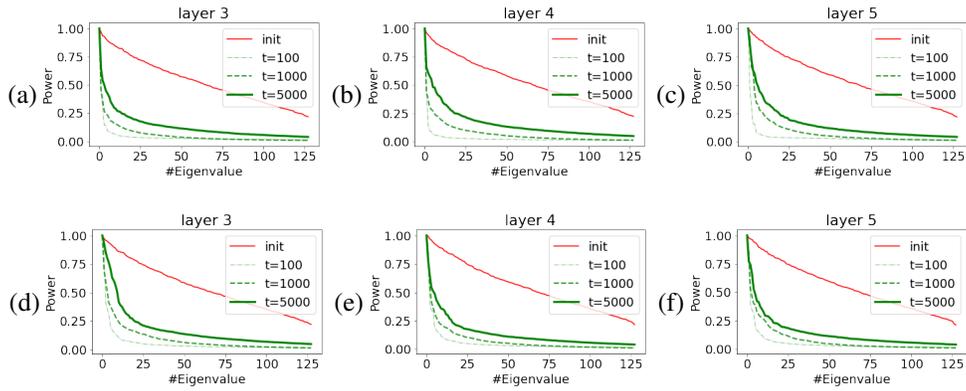


Figure 11: cifar-10: normalized spectrum at different stages of training. (a) - (c) Normalized spectrum of $\mathbf{W}_K \mathbf{W}_Q^\top$ at initialization and $\Delta \mathbf{W}_K \mathbf{W}_Q^\top$ during training for different attention heads at different layers. (d) - (e) equivalent figures for $\mathbf{W}_V \mathbf{W}_O^\top$.

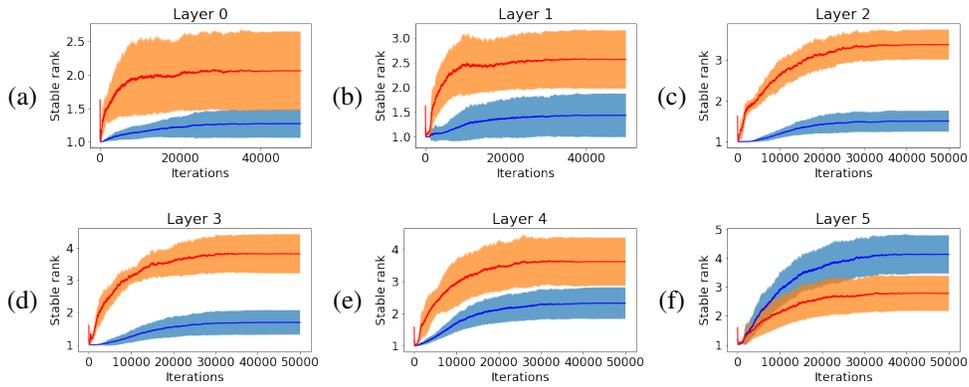


Figure 12: cifar-10: Stable rank of $\Delta \mathbf{W}_K \mathbf{W}_Q^\top$ (blue) and $\Delta \mathbf{W}_V \mathbf{W}_O^\top$ (red) throughout training. Mean and standard deviation (shaded area) are computed across 8 heads per attention layer.

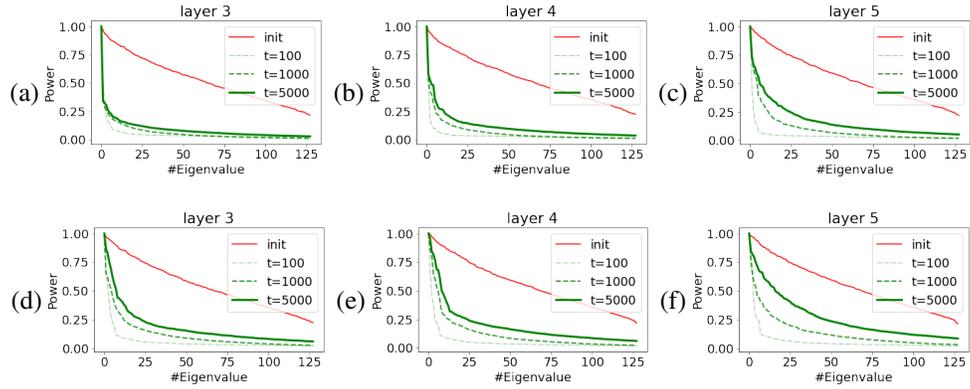


Figure 13: cifar-100: normalized spectrum at different stages of training. (a) - (c) Normalized spectrum of $\mathbf{W}_K \mathbf{W}_Q^\top$ at initialization and $\Delta \mathbf{W}_K \mathbf{W}_Q^\top$ during training for different attention heads at different layers. (d) - (e) equivalent figures for $\mathbf{W}_V \mathbf{W}_O^\top$.

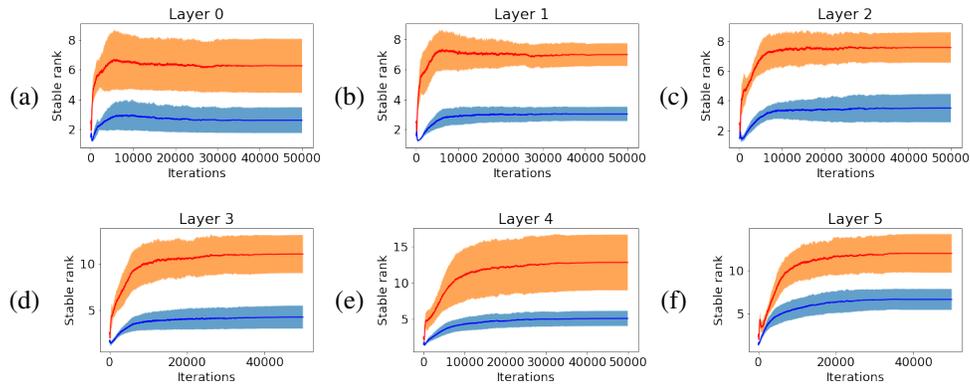


Figure 14: cifar-100: Stable rank of $\Delta \mathbf{W}_K \mathbf{W}_Q^\top$ (blue) and $\Delta \mathbf{W}_V \mathbf{W}_O^\top$ (red) throughout training. Mean and standard deviation (shaded area) are computed across 8 heads per attention layer.

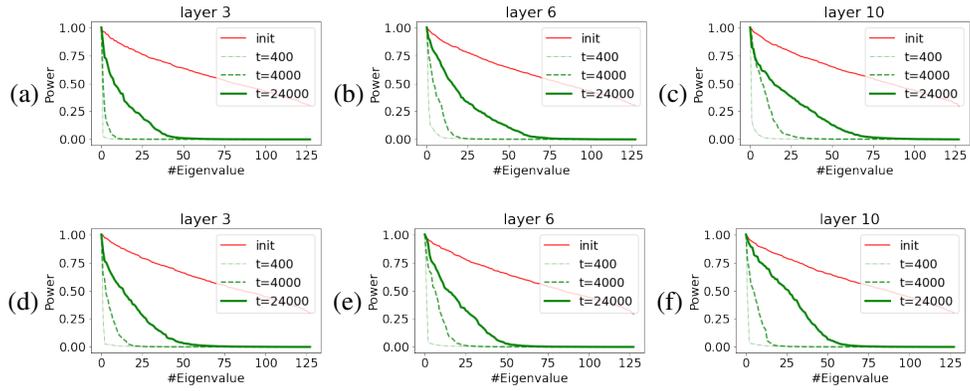


Figure 15: Imagenet: normalized spectrum at different stages of training. (a) - (c) Normalized spectrum of $\Delta W_K W_Q^\top$ at initialization and $\Delta W_K W_Q^\top$ during training for different attention heads at different layers. (d) - (e) equivalent figures for $\Delta W_V W_O^\top$.

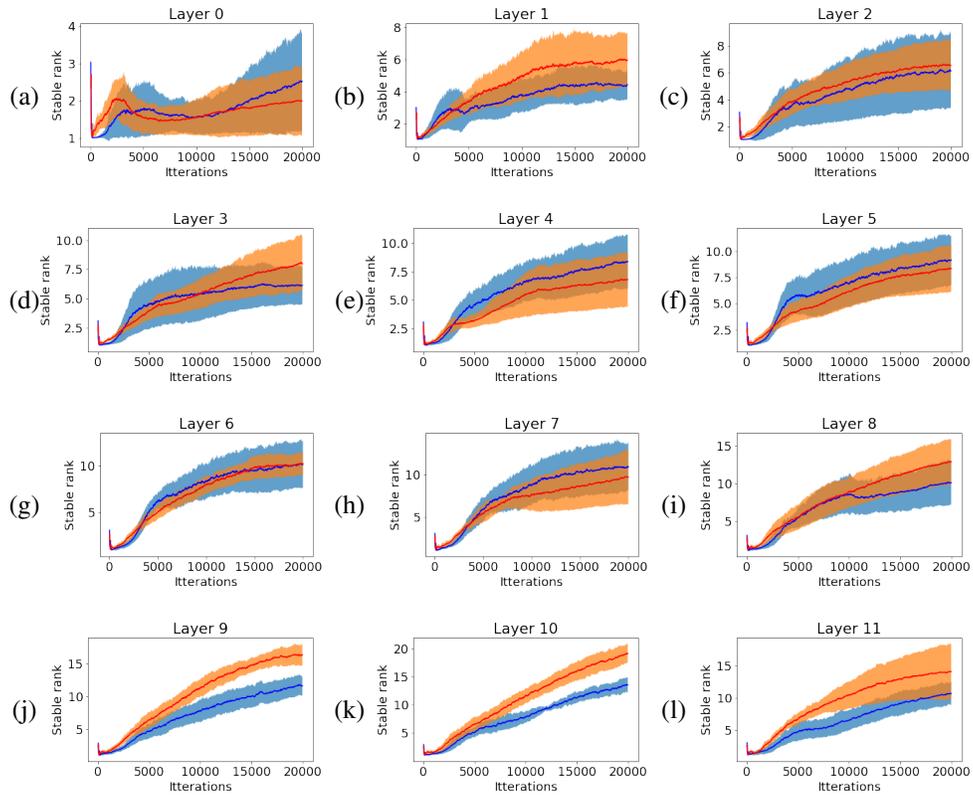


Figure 16: Imagenet: Stable rank of $\Delta W_K W_Q^\top$ (blue) and $\Delta W_V W_O^\top$ (red) throughout training. Mean and standard deviation (shaded area) are computed across 12 heads per attention layer.