IMPLICIT BIAS AND LOSS OF PLASTICITY IN MATRIX COMPLETION: DEPTH PROMOTES LOW-RANKNESS

Anonymous authorsPaper under double-blind review

ABSTRACT

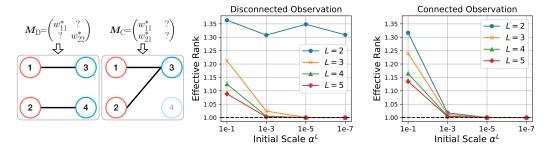
We study matrix completion via deep matrix factorization (a.k.a. deep linear neural networks) as a simplified testbed to examine how network depth influences training dynamics. Despite the simplicity and importance of the problem, prior theory largely focuses on shallow (depth-2) models and does not fully explain the implicit low-rank bias observed in deeper networks. We identify *coupled dynamics* as a key mechanism behind this bias and show that it intensifies with increasing depth. Focusing on gradient flow under diagonal observations, we prove: (a) networks of $depth \geq 3$ exhibit coupling unless initialized diagonally, and (b) convergence to rank-1 occurs if and only if the dynamics is coupled—resolving an open question by Menon (2024) for a family of initializations. We also revisit the loss of plasticity phenomenon in matrix completion (Kleinman et al., 2024), where pre-training on few observations and resuming with more degrades performance. We show that deep models avoid plasticity loss due to their low-rank bias, whereas depth-2 networks pre-trained under decoupled dynamics fail to converge to low-rank, even when resumed training (with additional data) satisfies the coupling condition shedding light on the mechanism behind this phenomenon.

1 Introduction

Overparameterized neural networks have the capacity to perfectly memorize the training data, even when they are given random labels (Zhang et al., 2017). Despite their large capacity, neural networks often generalize well to unseen data without any explicit regularization techniques, which challenges conventional statistical wisdom. Recent studies attribute this phenomenon to the implicit bias of neural networks, arguing that among the many possible global minima, first-order algorithms such as (stochastic) gradient descent favor solutions that generalize well (Neyshabur et al., 2014; 2017; Huh et al., 2021; Timor et al., 2023; Frei et al., 2023; Kou et al., 2023; Galanti et al., 2024; Jacot, 2022).

Matrix completion, a task with practical applications in areas like recommender systems and image restoration, provides a key framework for investigating these implicit biases, particularly the tendency towards low-rank solutions. While matrix completion can be viewed as a special case of the broader matrix sensing framework (Jin et al., 2023; Soltanolkotabi et al., 2023; Ma & Fattahi, 2023; Stöger & Soltanolkotabi, 2021; Li et al., 2018), which offers general tools for understanding recovery from limited data, specific challenges can emerge when applying these general theories directly. Notably, common theoretical assumptions prevalent in matrix sensing analyses, such as the Restricted Isometry Property (RIP) (Candes & Tao, 2005), often prove too stringent or may not adequately capture the nuances of many practical matrix completion tasks. For instance, even when completing the 2×2 matrix $M_{\rm C}$ (introduced in Figure 1a), which can successfully converge to a low-rank solution, the RIP condition cannot be satisfied. Therefore, researchers have investigated implicit bias phenomena specifically within matrix completion, without assuming the RIP condition (Menon, 2024; Bai et al., 2024; Razin & Cohen, 2020; Ma & Fattahi, 2024; Kim & Chung, 2023).

The goal of the matrix completion task is to recover a low-rank ground truth matrix W^* using only a subset of its entries. A common strategy for matrix completion involves matrix factorization, which can also be viewed as linear neural networks. These networks reparameterize the target matrix X as a product of factors, $X = W_L W_{L-1} \cdots W_1$, and train these factors W_i by minimizing the mean squared error on the observed entries via gradient descent. The observed entries constitute the training set, while the unobserved entries act as the test set.



(a) Bipartite graph of $M_{\rm D}$ & $M_{\rm C}$ (b) Effective rank trained w/ $M_{\rm D}$ (c) Effective rank trained w/ $M_{\rm C}$

Figure 1: (a) Examples of bipartite graphs corresponding to observation patterns of $M_{\rm D}$ (disconnected) and $M_{\rm C}$ (connected). (b-c) Training results showing effective rank (cf. Roy & Vetterli (2007)) for completing rank-1 matrices $M_{\rm D}$ and $M_{\rm C}$, respectively. The rank-1 ground truth matrices were generated via $uv^{\rm T}$, where $u,v\in\mathbb{R}^2$ with entries sampled i.i.d. from a standard normal distribution. We initialized each layer's entries by sampling from a Gaussian distribution with mean zero and standard deviation α , chosen to ensure the initial scale of the product matrix $W_{L:1}(0)$ is approximately invariant to depth L. Each result shows an average of 300 independent random trials.

The problem of predicting W^* is underdetermined, as infinitely many completions are possible. Nevertheless, both theory and experiments indicate that training even a simple two-layer factorization (L=2) with gradient descent, without explicit rank constraints, typically yields a low-rank solution under reasonable assumptions (Razin & Cohen, 2020; Bai et al., 2024; Ma & Fattahi, 2024).

A recent work by Bai et al. (2024) formalizes this phenomenon using the concept of *data connectivity*. They demonstrate that if the observed entries form a connected bipartite graph (meaning any observed entry can be reached from any other via shared rows or columns), a depth-2 factorization initialized at an infinitesimally small scale converges to a low-rank solution. Conversely, the network may converge to a higher-rank matrix if the observations are disconnected (see Definition 1 and Figure 1a).

However, the situation changes significantly for deeper ($L \geq 3$) networks, as empirically demonstrated in Figure 1. Consider the task of completing the 2×2 matrix

$$\boldsymbol{M}_{\mathrm{D}} = \begin{pmatrix} w_{11}^* & ? \\ ? & w_{22}^* \end{pmatrix} \tag{1}$$

where only the diagonal entries are observed. This observation pattern forms a disconnected graph as illustrated in Figure 1a. Consistent with the theory for disconnected graphs, L=2 models fail to find a low-rank solution, empirically converging to rank-2 regardless of initialization scale. In contrast, deeper models ($L\geq 3$) with small initialization tend to converge to a rank-1 solution, as shown in Figure 1b. This specific example highlights that the implicit low-rank bias appears to be strengthened by depth, in a way that cannot be explained solely by the data connectivity framework developed for L=2 models. Furthermore, considering connected cases as well, Figure 1c demonstrates that this strong low-rank bias is generally robust, tending to strengthen further as depth increases.

However, a theoretical understanding of this depth-induced bias remains elusive, largely due to the complex, coupled dynamics during training. While Arora et al. (2019) offer insights, their claim that the gap between two arbitrary singular values widens with depth is not fully formal. It stems largely from their analysis assuming stabilized singular vectors, which limits its scope. Indeed, Menon (2024) notes that even for a simple case like (1) with $w_{11}^* = w_{22}^* = 1$, proving that gradient descent with a deep factorization converges to a low-rank solution is still an open problem. Motivated by this gap in understanding, we theoretically analyze such settings, including the example (1).

Investigating the implicit low-rank bias in matrix completion can also shed light on the phenomenon of "loss of plasticity", a challenge widely observed in general neural network training (Shin et al., 2024; Ash & Adams, 2020; Achille et al., 2018; Berariu et al., 2021). The term loss of plasticity describes the tendency of neural networks, particularly after initial training, to lose their adaptability to new information, hindering their generalization capabilities. A recent work by Kleinman et al. (2024) empirically reports this phenomenon even in matrix completion. They observe that models

trained with insufficient data often yield high-rank solutions. If these models then warm-start using augmented data, they frequently struggle to achieve low-rank solutions. To provide a theoretical explanation for why this loss of plasticity occurs, this paper elucidates the phenomenon.

To summarize, here are the main research questions that we address throughout the paper:

- What is the fundamental difference between deep $(L \ge 3)$ and shallow (L = 2) factorizations regarding their implicit low-rank bias, particularly for disconnected observations?
- Can we theoretically establish that deeper models (i.e., with larger L ≥ 3) exhibit a stronger implicit bias toward low-rank solutions?
- What is the underlying cause of the loss of plasticity phenomenon, and how does depth interplay with it?

In Section 3.1, we begin by examining the depth-2 case to elucidate the key mechanism of connectivity. We find that *coupled training dynamics* induces a low-rank bias, a phenomenon generalizable to deeper networks. Section 3.2 further investigates this for all $L \geq 2$ using the diagonal observation case. Our analysis reveals that, for deep models, this bias distinctively promotes low-rank solutions compared to depth-2 models, strengthening with depth. Finally, Section 4 explores the loss of plasticity phenomenon in matrix completion. We observe that deep models typically avoid this phenomenon due to their low-rank bias. In contrast, we empirically observe and prove that depth-2 networks pre-trained with limited observations (yielding decoupled dynamics) and subsequently trained with augmented observations (yielding coupled dynamics) fail to find a low-rank solution.

2 Problem Setting

We consider the problem of estimating a ground truth matrix $\boldsymbol{W}^* \in \mathbb{R}^{d \times d}$ based on observations of its entries $\{w_{ij}^*\}_{(i,j)\in\Omega}$, where $\Omega \subseteq [d] \times [d]$ is the set of observed indices. We model the estimate as a linear network $\boldsymbol{W}_{L:1} \triangleq \boldsymbol{W}_L \boldsymbol{W}_{L-1} \cdots \boldsymbol{W}_1$, where $\boldsymbol{W}_l \in \mathbb{R}^{d_l \times d_{l-1}}$ with $d_0 = d_L = d$. We denote the (i,j)-th entry of the matrix $\boldsymbol{W}_{L:1}$ as w_{ij} . The factor matrices $\{\boldsymbol{W}_l\}_{l=1}^L$ are trained by minimizing an objective function ϕ , defined as the mean squared error ℓ over the observed entries in Ω :

$$\phi(\mathbf{W}_1, \dots, \mathbf{W}_L; \Omega) \triangleq \ell(\mathbf{W}_{L:1}; \Omega) = \frac{1}{2} \sum_{(i,j) \in \Omega} \left(w_{ij} - w_{ij}^* \right)^2.$$
 (2)

We study the overparameterized regime where the intermediate dimensions satisfy $d_l \geq d$ for all $l \in [L-1]$, imposing no explicit rank constraints on the product model $\mathbf{W}_{L:1}$. Consistent with prior works, our analysis focuses on *gradient flow* dynamics (gradient descent with an infinitesimal step size) for a given objective function ϕ . The dynamics for each layer $\mathbf{W}_l(t)$ evolve according to:

$$\dot{\mathbf{W}}_{l}(t) \triangleq \frac{d}{dt}\mathbf{W}_{l}(t) = -\frac{\partial}{\partial \mathbf{W}_{l}(t)}\phi(\mathbf{W}_{1}(t), \mathbf{W}_{2}(t), \dots, \mathbf{W}_{L}(t); \Omega), \quad l \in [L], \ t \geq 0.$$
 (3)

For depth-2 networks (L=2), the product of factor matrices $\boldsymbol{A} \in \mathbb{R}^{d \times d_1}$ (representing \boldsymbol{W}_2) and $\boldsymbol{B} \in \mathbb{R}^{d_1 \times d}$ (representing \boldsymbol{W}_1), we denote $\boldsymbol{W}_{\boldsymbol{A},\boldsymbol{B}} \triangleq \boldsymbol{A}\boldsymbol{B}$.

Bai et al. (2024) introduce the concept of data connectivity for an incomplete matrix M. Connectivity is characterized by its set of observed indices $\Omega \subseteq [d] \times [d]$ and the corresponding observation matrix P (where $P_{ij} = 1$ if $(i, j) \in \Omega$, and 0 otherwise). The formal definition is as follows:

Definition 1 (Connectivity from Bai et al. (2024)). An incomplete matrix M is connected if the bipartite graph \mathcal{G}_M , constructed from its observation matrix P using the adjacency matrix $\begin{bmatrix} 0 & P^\top \\ P & 0 \end{bmatrix}$, is connected after removing isolated vertices. Otherwise, M is disconnected.

3 IMPLICIT BIAS OF DEPTH INDUCED BY COUPLED TRAINING DYNAMICS

In this section, we extend the connectivity argument of Bai et al. (2024) to general depth factorizations. We first demonstrate how the *coupling of training dynamics* serves as the key mechanism explaining data connectivity's role in depth-2 models, through the completion of two previously introduced 2×2 matrices, $M_{\rm D}$ and $M_{\rm C}$, as illustrative examples. Building on the insights derived from these depth-2 model analyses, we hypothesize that deep networks exhibit an intrinsic low-rank bias because they maintain a high degree of coupled training dynamics, irrespective of observation patterns. This hypothesis is further corroborated by the diagonal observation results presented in Section 3.2.

3.1 WARM-UP: COUPLED DYNAMICS VS. DECOUPLED DYNAMICS IN DEPTH-2 NETWORKS

We focus on the simple 2×2 matrix completion of M_D and M_C , using depth-2 models $W_{A,B}(t) = A(t)B(t)$. For brevity, let $a_i(t) \in \mathbb{R}^{d_1}$ be the transpose of the *i*-th row of A(t), and let $b_j(t) \in \mathbb{R}^{d_1}$ be the *j*-th column of B(t). Our aim is to see how training dynamics affect the *alignment* of the rows of A(t) or the columns of B(t), as such alignment leads to a rank-1 product matrix $W_{A,B}(t)$.

Decoupled Dynamics. In the M_D case (disconnected observations w_{11}^*, w_{22}^*), the gradient flow using the objective defined in (2), results in independent dynamics for the pairs (a_1, b_1) and (a_2, b_2) :

$$\dot{\boldsymbol{a}}_i(t) = (w_{ii}^* - \boldsymbol{a}_i(t)^\top \boldsymbol{b}_i(t)) \boldsymbol{b}_i(t), \quad \dot{\boldsymbol{b}}_i(t) = (w_{ii}^* - \boldsymbol{a}_i(t)^\top \boldsymbol{b}_i(t)) \boldsymbol{a}_i(t) \quad \text{for } i = 1, 2.$$

Note that while the dynamics couple $a_1(t)$ with $b_1(t)$ and $a_2(t)$ with $b_2(t)$ within each pair, the two pairs (a_1,b_1) and (a_2,b_2) are decoupled. This decoupling means the overall system's dynamics separate into two independent systems. Consequently, there is no compelling reason to align vectors from different pairs, typically leading to high-rank solutions with generic initializations (Figure 1b). Indeed, we can obtain closed-form solutions solely dependent on initialization (see Proposition 4.1). For instance, with $A(0) = B(0) = \alpha I_2$, we have $W_{A,B}(\infty) = \operatorname{diag}(w_{11}^*, w_{22}^*)$, a rank-2 solution.

Coupled Dynamics. In contrast, for the M_C case (connected observations w_{11}^*, w_{21}^*), the gradient flow on the objective (2) yields coupled dynamics that do not decompose into independent pairs:

$$\dot{a}_{1}(t) = (w_{11}^{*} - a_{1}(t)^{\top} b_{1}(t)) b_{1}(t), \quad \dot{a}_{2}(t) = (w_{21}^{*} - a_{2}(t)^{\top} b_{1}(t)) b_{1}(t),
\dot{b}_{1}(t) = (w_{11}^{*} - a_{1}(t)^{\top} b_{1}(t)) a_{1}(t) + (w_{21}^{*} - a_{2}(t)^{\top} b_{1}(t)) a_{2}(t).$$
(4)

An important observation from (4) is that A(0) = 0 ensures rank-1 $W_{A,B}(t)$ due to persistent alignment of $a_1(t)$, $a_2(t)$ and $b_1(t)$. Although non-zero initialization leads to more complex behavior arising from coupled training dynamics, the following theorem demonstrates that sufficiently small initial norms in A(0) also result in the alignment of $a_1(t)$ and $a_2(t)$ with $b_1(t)$.

Theorem 3.1. For the product model $W_{A,B}(t) = A(t)B(t) \in \mathbb{R}^{2\times 2}$, we consider the gradient flow dynamics (4), where the observations are $w_{11}^*(\neq 0)$ and $w_{21}^*(\neq 0)$. We assume convergence to the zero-loss solution (i.e., $w_{11}(\infty) = w_{11}^*, w_{21}(\infty) = w_{21}^*$). Defining $\mathbf{u}^* = \frac{\mathbf{b}_1(\infty)}{\|\mathbf{b}_1(\infty)\|_2}$ and the orthogonal component $\mathbf{a}_{i\perp}(\infty) = \mathbf{a}_i(\infty) - (\mathbf{a}_i(\infty)^{\top}\mathbf{u}^*)\mathbf{u}^*$, we have:

$$\frac{\|\boldsymbol{a}_{i\perp}(\infty)\|_{2}^{2}}{\|\boldsymbol{a}_{i}(\infty)\|_{2}^{2}} \leq \frac{\|\boldsymbol{A}(0)\|_{F}^{2} \left(\sqrt{\|\boldsymbol{b}_{1}(0)\|_{2}^{4} + 4{w_{11}^{*}}^{2} + 4{w_{21}^{*}}^{2}} + \|\boldsymbol{b}_{1}(0)\|_{2}^{2}\right)}{2{w_{i1}^{*}}^{2}}, \ \textit{for } i = 1, 2.$$

The theorem shows that small initial norms for A(0) lead to the alignment of $a_1(\infty)$ and $a_2(\infty)$ with $b_1(\infty)$, implying a near rank-1 product matrix $W_{A,B}(\infty)$. This suggests that for depth-2 networks, coupled training dynamics (resulting from connected observations) facilitate the emergence of low-rank solutions under such small initialization, in contrast to the decoupled dynamics of disconnected observations, where no such bias exists regardless of initialization scale. This connection between observation connectivity and the coupling of training dynamics in depth-2 models motivates our investigation into how coupled dynamics manifest and induce low-rank bias in deeper networks, irrespective of connectivity patterns, as explored in the subsequent sections.

Remark. Analyzing these dynamics is challenging because the time evolutions of a_1 , a_2 , and b_1 are mutually dependent. We note that Theorem 3.1 is not a direct corollary of Theorem 3 in Bai et al. (2024). We explicitly characterize the degree of misalignment as a function of the initialization scale, unlike their assumption of an infinitesimal initialization scale with additional conditions.

3.2 COUPLED DYNAMICS IN DEEP NETWORKS INDUCE IMPLICIT BIAS TOWARDS LOW RANK

Section 3.1 illustrated the importance of coupled training dynamics, driven by data connectivity, for achieving low-rank solutions in simple two-layer factorizations (L=2). Building on this understanding, we now extend our analysis to deep networks ($L\geq 3$). For illustrative purposes, consider a depth-3 network $W_{3:1}$. An arbitrary observed entry w_{ij} from this matrix is given by:

$$w_{ij} = \sum_{k=1}^{d_2} \sum_{l=1}^{d_1} (\mathbf{W}_3)_{ik} (\mathbf{W}_2)_{kl} (\mathbf{W}_1)_{lj}.$$
 (5)

Crucially, because all elements of the intermediate matrix W_2 contribute to the computation of w_{ij} regardless of (i,j), gradients of different observed entries will propagate through and update these shared elements in W_2 . This inherently couples their training dynamics, a structural feature distinct from the depth-2 case, where coupling is primarily determined by the observation pattern. Such inherent coupling, in turn, implies a potential intrinsic bias towards low-rank solutions for deep models. To formalize this notion, we introduce the following definition of coupled dynamics.

Definition 2 (Coupled/Decoupled Dynamics). Consider the matrix completion setup with the model $W_{L:1}(t) = W_L(t) \cdots W_1(t) \in \mathbb{R}^{d \times d}$. Let $\theta(t)$ be the vector of all trainable parameters evolving according to the gradient flow dynamics (defined in (3)). The gradient flow dynamics are **decoupled** if there exists a partition of Ω into non-empty, disjoint subsets $\Omega_1, \ldots, \Omega_K$ ($K \geq 2$) such that $\bigcup_{k=1}^K \Omega_k = \Omega$ and the following condition holds for any $(i,j) \in \Omega_k$ and $(p,q) \in \Omega_l$ with $k \neq l$:

$$\langle \nabla_{\boldsymbol{\theta}} w_{ij}(t), \nabla_{\boldsymbol{\theta}} w_{pq}(t) \rangle = 0, \quad \forall t \ge 0.$$
 (6)

The gradient flow dynamics are coupled if they are not decoupled.

For depth-2 matrices, it is straightforward to verify that coupled and decoupled dynamics typically correspond to connected and disconnected graphs, respectively, based on Definitions 1 and 2. For depth ≥ 3 matrices, any initialization with an absolutely continuous distribution (e.g., Gaussian, uniform) yields gradient flow dynamics that are coupled with probability one (see Proposition B.1 in Appendix B), irrespective of the observation pattern. However, special cases exist where training dynamics are decoupled even for $L \geq 3$. Refer to Appendix B for further discussion.

3.2.1 IMPLICIT BIAS OF DEPTH UNDER DIAGONAL OBSERVATIONS

To gain deeper theoretical insight into how coupled dynamics induce low-rank bias as depth increases, we further investigate the diagonal observation setting. As highlighted in the 2×2 example (cf. Figure 1b), this setting reveals a stark difference between shallow and deep networks despite being a disconnected observation pattern. To investigate this further, we now turn to the general $d \times d$ case.

Specifically, we consider a $d \times d$ ground truth matrix \boldsymbol{W}^* with positive and identical diagonal observations $w^* \triangleq w_{11}^* = \cdots = w_{dd}^* > 0$ where $\Omega_{\mathrm{diag}}^{(d)} \triangleq \{(i,i) \mid i \in [d]\}$. We factorize the model with depth-L: $\boldsymbol{W}_{L:1}(t) = \boldsymbol{W}_L(t) \boldsymbol{W}_{L-1}(t) \cdots \boldsymbol{W}_1(t)$ where $\boldsymbol{W}_l \in \mathbb{R}^{d \times d}$ for all $l \in [L]$.

To investigate how dynamic coupling affects the low-rank bias, we consider a family of initializations where, for parameters $\alpha > 0$ and m > 1, each factor matrix $\mathbf{W}_l(0)$ is initialized as follows:

$$\mathbf{W}_{l}(0) = \begin{pmatrix} \alpha & \alpha/m & \cdots & \alpha/m \\ \alpha/m & \alpha & \cdots & \alpha/m \\ \vdots & \vdots & \ddots & \vdots \\ \alpha/m & \alpha/m & \cdots & \alpha \end{pmatrix} \in \mathbb{R}^{d \times d}, \quad \forall l \in [L].$$
 (7)

Using this initialization scheme with diagonal observations, the following proposition specifies how parameters m and network depth L determine if training dynamics are coupled or decoupled:

Proposition 3.2. Consider a depth-L model, where each factor $W_l(0) \in \mathbb{R}^{d \times d}$ is initialized with (7) trained with diagonal observations, $\Omega_{\text{diag}}^{(d)}$. Then, according to Definition 2, the following hold:

- For depth L=2, the training dynamics are **decoupled** for all m>1.
- For depth $L \geq 3$:
 - The training dynamics are **coupled** if $1 < m < \infty$.
 - The training dynamics are **decoupled** if $m = \infty$ (i.e., initialization with αI_d).

By Proposition D.1 in Appendix D, the loss decays exponentially to zero under the gradient flow dynamics (3). Building on this zero-loss convergence, our objective is to determine the rank of solutions found by gradient flow depending on the coupling of dynamics. The theorem below presents an equation of each singular value of the converged matrix $W_{L:1}(\infty)$, for all $L \ge 2$.

Theorem 3.3. Consider the product matrix $W_{L:1}$, whose factor matrices $W_l \in \mathbb{R}^{d \times d}$ are initialized according to (7). Under the gradient flow dynamics (3), we have $\ell(W_{L:1}(\infty); \Omega_{\mathrm{diag}}^{(d)}) = 0$ (Proposition D.l, Appendix D). Let $\sigma_1 \geq \cdots \geq \sigma_d \geq 0$ denote the singular values of the converged matrix $W_{L:1}(\infty)$. Then, for all parameter values $\alpha > 0$, m > 1, $d \geq 2$, and $L \geq 2$, the following holds:

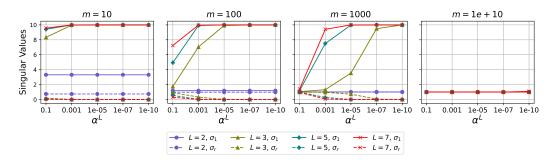


Figure 2: Singular values σ_i of $W_{L:1}(\infty)$ (numerically obtained from Theorem 3.3) against initialization scale α^L , for the diagonal observation task. Solid lines represent the largest singular value σ_1 ; dashed lines denote the other (identical) singular values σ_r for $r \geq 2$. For finite m, these results illustrate that both greater depth L and a smaller initial scale α enhance the low-rank bias, in contrast to the L=2 case. Conversely, a very large m (e.g., $m=10^{10}$), approximating an αI_d (rank-d) initialization, leads to decoupled dynamics and a full-rank solution, independent of both L and α .

- If L=2 (decoupled dynamics): The singular values are explicitly given by

$$\sigma_1 = \frac{w^*(m+d-1)^2}{m^2+d-1}, \quad \sigma_r = \frac{w^*(m-1)^2}{m^2+d-1} \quad \text{for } r = 2, \dots, d.$$

- If $L \ge 3$ and $1 < m < \infty$ (coupled dynamics): The singular values satisfy the implicit equations:

$$(\sigma_1)^{\frac{2-L}{L}} - \left(\frac{w^*d - \sigma_1}{d - 1}\right)^{\frac{2-L}{L}} = C_{\alpha, m, L, d},\tag{8}$$

$$(w^*d - (d-1)\sigma_r)^{\frac{2-L}{L}} - (\sigma_r)^{\frac{2-L}{L}} = C_{\alpha,m,L,d}, \quad \text{for } r = 2, \dots, d,$$
(9)

where $C_{\alpha,m,L,d} \triangleq \left(\frac{\alpha}{m}\right)^{2-L} \left((m+d-1)^{2-L}-(m-1)^{2-L}\right)$.

- If $L \ge 3$ and $m = \infty$ (decoupled dynamics): The singular values converge to:

$$\sigma_i = w^*, \text{ for } i = 1, 2, \dots d.$$

The proof of the theorem is provided in Appendix D.3. The theorem details the converged singular values of $W_{L:1}(\infty)$ for our initialization scheme (7). Crucially, it reveals distinct outcomes based on the nature of the training dynamics. For decoupled dynamics—specifically, when L=2 (for sufficiently large m>1), or when $L\geq 3$ and $m=\infty$ —all singular values approach w^* and are independent of the scale α . This implies convergence to a full-rank solution. In contrast, for coupled dynamics ($L\geq 3$ with finite m), the outcome becomes α -dependent. The analytical intractability of of the governing implicit equations in this coupled regime motivates a numerical study.

To numerically investigate this, we solve the implicit equations (8) and (9) that determine singular values σ_i for the coupled $L \geq 3$, finite m case. Setting $w^* = 1$ and d = 10, we examine how network depth (L) and initialization parameters (α, m) influence the singular value distribution. The results (Figure 2) confirm that these coupled dynamics in models with $L \geq 3$ and finite m indeed induce a low-rank bias, contrasting with the full-rank outcomes of the decoupled cases. Moreover, this bias becomes more pronounced as L increases, evidenced by a wider gap between σ_1 and σ_r for $r \geq 2$.

Additional numerical evidences are provided in Figures 5–7 (Appendix C.1). Moreover, Figure 8 in Appendix C.1 shows that these numerical results agree with the outcomes of a gradient descent with a sufficiently small learning rate. We further train practical neural networks to examine whether increased depth indeed leads to a low-rank bias. The results shown in Figures 10–13 (Appendix C.1.1) indicate that as depth increases (e.g., ResNet-18 to 101 and VGG-11 to 19), the average effective rank decreases, highlighting the emergence of low-rank bias in practical neural networks.

Remark. Our analysis of low-rank bias for a specific family of deterministic initializations resolves the challenging open problem (1) highlighted in Section 14.1 of Menon (2024). Figure 9 in Appendix C.1 further demonstrates that our proposed deterministic initialization exhibits qualitative trends similar to Gaussian initialization. We therefore argue that our results provide foundational insights into low-rank bias applicable to more general random initializations.

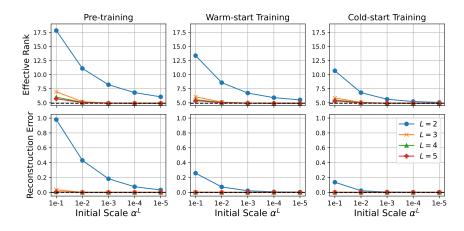


Figure 3: Experiments use a 100×100 rank-5 ground-truth matrix. pre-training utilizes 2000 randomly sampled entries ($\Omega_{\rm pre}$; $|\Omega_{\rm pre}| = 2000$), while post-training adds 1000 more, forming $\Omega_{\rm post}$ ($\Omega_{\rm pre} \subset \Omega_{\rm post}$; $|\Omega_{\rm post}| = 3000$). The top row of panels displays effective rank, and the bottom row shows reconstruction error, both measured at convergence. The leftmost panels depict training on $\Omega_{\rm pre}$, and the rightmost on $\Omega_{\rm post}$, both starting from random Gaussian initialization. The middle panels show warm-start training on $\Omega_{\rm post}$, initialized from converged pre-trained models with $\Omega_{\rm pre}$.

4 Understanding Loss of Plasticity in Depth-2 Matrix Completion

Studying the inherent tendency towards low-rank solutions in matrix completion can offer further insights into the loss of plasticity phenomenon. Kleinman et al. (2024) report the emergence of this phenomenon in matrix completion: models pre-trained on limited observations struggle to adapt when training continues on augmented observations. Notably, they observe that loss of plasticity is further intensified with increasing network depth, a conclusion they reached by measuring a "relative reconstruction loss" when compared to models trained from scratch on the augmented dataset.

However, our findings (Figure 3) offer a more nuanced perspective. We observed that even when pre-trained with a sparser set of observations, deeper models increasingly favor low-rank solutions as their depth increases. This aligns with our argument (Section 3.2) that they inherently achieve low-rank solutions even from limited, disconnected initial data. Consequently, for these deeper models, further training on augmented data (the post-training stage) does not lead to noticeably higher rank compared to training equivalent models from scratch on the augmented observations. Therefore, while their performance might exhibit a relative degradation compared to models trained from scratch, their absolute solution quality can still surpass that of shallower models. Based on our observations, we conclude that the low-rank bias of deep models helps them mitigate the loss of plasticity, while the phenomenon is more pronounced in depth-2 models. To theoretically understand the underlying cause of this phenomenon itself, we henceforth focus our analysis on depth-2 models.

In Section 4.1, we study pre-training on diagonal-only observations, i.e., the disconnected index set $\Omega_{\mathrm{diag}}^{(d)}$. We then consider post-training on 2×2 (Section 4.2) and $d \times d$ (Section 4.3) matrices. For the 2×2 case, we set $\Omega_{\mathrm{pre}}^{(2)} \triangleq \Omega_{\mathrm{diag}}^{(2)}$ and obtain the post-training set $\Omega_{\mathrm{post}}^{(2)}$ by adding a single off-diagonal entry to ensure connectivity. Likewise, for the $d \times d$ case, $\Omega_{\mathrm{pre}}^{(d)} \triangleq \Omega_{\mathrm{diag}}^{(d)}$, and $\Omega_{\mathrm{post}}^{(d)}$ is formed by adding additional (off-diagonal) observations; see Section 4.3 for details.

4.1 PRE-TRAINING WITH DIAGONAL OBSERVATIONS

To clearly observe loss of plasticity in a setting consistent with Section 3.2, we pre-train using only diagonal entries, yielding a disconnected pattern. We consider decoupled-to-coupled scenarios, where additional data is introduced to induce coupled training dynamics. For depth-2 models, they correspond to a disconnected-to-connected observation pattern. For the pre-training, closed-form solutions that depend *solely* on the network's initialization can be found in the following proposition:

 Proposition 4.1. Consider a ground truth matrix $\mathbf{W}^* \in \mathbb{R}^{d \times d}$ with diagonal observations $\Omega_{\mathrm{diag}}^{(d)}$. The model is factorized as $\mathbf{W}_{\mathbf{A},\mathbf{B}}(t) = \mathbf{A}(t)\mathbf{B}(t)$, where $\mathbf{A}(t),\mathbf{B}(t) \in \mathbb{R}^{d \times d}$. For each observation $(i,i) \in \Omega_{\mathrm{diag}}^{(d)}$, define the constants P_i and Q_i based on the initial values:

$$P_i \triangleq \sum_{k=1}^d a_{ik}(0)b_{ki}(0)$$
 and $Q_i \triangleq \sum_{k=1}^d (a_{ik}(0)^2 + b_{ki}(0)^2)$.

Furthermore, for each diagonal observation, let the parameter \bar{r}_i be determined from the ground truth

entry
$$w_{ii}^*$$
 and the constants defined above, $\bar{r}_i \triangleq \frac{1}{2} \log \left(\frac{P_i + \frac{Q_i}{2}}{w_{ii}^* + \sqrt{w_{ii}^*^2 - P_i^2 + \left(\frac{Q_i}{2}\right)^2}} \right)$. Then, assuming

convergence to a zero-loss solution of the loss $\ell(W_{A,B}; \Omega_{\mathrm{diag}}^{(d)})$, any entry $a_{pq}(\infty)$ of the converged matrix $A(\infty)$ and any entry $b_{pq}(\infty)$ of the converged matrix $B(\infty)$ (for any $p, q \in [d]$) are given by:

$$a_{pq}(\infty) = a_{pq}(0)\cosh(\bar{r}_p) - b_{qp}(0)\sinh(\bar{r}_p),$$

$$b_{pq}(\infty) = b_{pq}(0)\cosh(\bar{r}_q) - a_{qp}(0)\sinh(\bar{r}_q).$$

Remark. The proposition covers *arbitrary* initializations with *distinct* w_{ii}^* , which goes beyond Theorem 3.3 in the L=2 setting. While the above analysis focuses on diagonal observation cases, it can be generalized to any fully disconnected case (i.e., a single observation per row and column). This yields distinct solutions for various types of observation sets, as detailed in Appendix E.1.

We analyze the scenario where training resumes from a state obtained through pre-training. Let the pre-training phase conclude at a sufficiently large timestep T_1 . For simplicity, we assume that the solution $W_{A,B}(T_1)$ has perfectly converged with respect to the pre-training objective, neglecting any residual error due to the finite duration of this phase. Our subsequent analysis demonstrates that, starting from $W_{A,B}(T_1)$, the model $W_{A,B}(t)$ cannot converge to a low-rank solution.

4.2 Post-training: 2 by 2 Matrix Example

We aim to analyze scenarios where training is resumed under coupled dynamics, building upon solutions obtained from an initial decoupled pre-training phase (Proposition 4.1). To this end, we first define the specific pre-training setup for an illustrative 2×2 case: We observe diagonal entries $(\Omega_{\text{pre}}^{(2)})$, which are identical and positive, i.e., $w^* \triangleq w_{11}^* = w_{22}^* > 0$. To make loss of plasticity particularly pronounced during the pre-training, we initialize the model with αI_2 (for $\alpha > 0$), which is the $m = \infty$ setting of our initialization scheme in (7). Then, from Proposition 4.1, it follows that:

$$\mathbf{A}(T_1) = \mathbf{B}(T_1) = \begin{pmatrix} \sqrt{w^*} & 0\\ 0 & \sqrt{w^*} \end{pmatrix}. \tag{10}$$

For the subsequent post-training phase, an additional off-diagonal observation is introduced to establish connectivity. Without loss of generality, we assume $w_{12}^*>0$ is revealed, while the diagonal entries w_{11}^* and w_{22}^* from the pre-training phase remain observed. Thus, the updated set of observed entries becomes $\Omega^{(2)}_{\rm post}=\{(1,1),(1,2),(2,2)\}$. The ground-truth matrix is assumed to be rank-1, ensuring the setting is non-trivial, and the task is thus to predict the remaining entry $w_{21}^*=w^{*2}/w_{12}^*>0$. The following theorem, however, reveals a contrasting outcome for this entry.

Theorem 4.2. Let $A(T_1)$, $B(T_1)$ be the factor matrices obtained from the pre-training phase, as specified by (10). Then, running gradient flow during the subsequent post-training phase (for $t \ge T_1$), starting from $A(T_1)$ and $B(T_1)$, results in exponential decay of the loss:

$$\ell(\boldsymbol{W}_{A,B}(t); \Omega_{\text{post}}^{(2)}) \le \frac{1}{2} w_{12}^{*2} e^{-2w^{*}(t-T_{1})}.$$

Consequently, a lower bound for the stable rank of the converged matrix $W_{A,B}(\infty)$ is given by:

$$\frac{\|\boldsymbol{W_{A,B}(\infty)}\|_F^2}{\|\boldsymbol{W_{A,B}(\infty)}\|_2^2} \ge 1 + \exp\left(-8\frac{w_{12}^*}{w^*}\right).$$

Furthermore, for all $t > T_1$, $w_{21}(t)$ of the evolving matrix $\mathbf{W}_{A,B}(t)$ satisfies $w_{21}(t) < 0$.

The theorem indicates that the loss decreases exponentially fast, particularly when starting from large-norm solutions (at a rate governed by w^*). Therefore, since the model converged to high-rank solutions during pre-training, its singular values remain largely unchanged from this initial state, as long as w_{12}^* has a small magnitude compared to w^* . Furthermore, the unobserved entry $w_{21}(t)$ converges to a negative value, which contradicts the positive w_{21}^* expected for the true rank-1 solution.

4.3 Post-training: d by d Matrix under Lazy Training Regime

We attribute Theorem 4.2 primarily to the model's "lazy training" (Chizat et al., 2019) as large-norm initializations lead to faster loss decay, causing the model to converge to a nearby global minimum that may not be a low-rank solution. Drawing on this concept, we extend the preceding analysis of loss of plasticity to the more general case of $d \times d$ ground-truth matrices. The following theorem states that when the model is initialized with a sufficiently small loss, resulting from warm-starting that perfectly fits all previously observed data, the model exhibits lazy training. This, in turn, prevents further learning that would reduce the rank and instead steers the model towards a nearby minimum.

Theorem 4.3. For factor matrices $A, B \in \mathbb{R}^{d \times d}$, suppose A and B are balanced at t = 0, i.e., $A(0)^{\top}A(0) = B(0)B(0)^{\top}$. Let f(A, B) be the function that maps (A, B) to the vector of model predictions for a given set of observed entries $\Omega_{\text{post}}^{(d)}$. We then define σ_{max} and σ_{min} as the maximum and minimum singular values, respectively, of the Jacobian of the function f evaluated at the pretrained state (at $t = T_1$). If the loss at time T_1 satisfies $\ell\left(W_{A,B}(T_1); \Omega_{\text{post}}^{(d)}\right) \leq \frac{\sigma_{\text{min}}^6}{1152d\sigma_{\text{max}}^2}$, this results in exponential decay of the loss:

$$\ell\left(\boldsymbol{W_{A,B}}(t);\Omega_{\mathrm{post}}^{(d)}\right) \leq \ell\left(\boldsymbol{W_{A,B}}(T_1);\Omega_{\mathrm{post}}^{(d)}\right) \exp\left(-\frac{1}{2}\sigma_{\min}^2(t-T_1)\right).$$

Consequently, the stable rank of A(t) (which is equal to that of B(t)) remains bounded below by

$$\frac{\|\boldsymbol{A}(t)\|_F^2}{\|\boldsymbol{A}(t)\|_2^2} \ge \left(\frac{\|\boldsymbol{A}(T_1)\|_F - \frac{\sigma_{\min}}{4\sqrt{2d}}}{\|\boldsymbol{A}(T_1)\|_2 + \frac{\sigma_{\min}}{4\sqrt{2d}}}\right)^2.$$

The theorem states that if a model has little remaining to learn (achieved via pre-training), it undergoes lazy training regime. In this regime, the loss converges rapidly, while its stable rank remains largely unchanged from the initial state. Thus, once a model has converged to a high-rank state, it struggles to recover a low-rank structure even when new observations are introduced to form connectivity. The proof of Theorem 4.3 is provided in Appendix E.3.

Example. As an illustrative example, consider a rank-1 ground-truth matrix $W^* \in \mathbb{R}^{d \times d}$,

$$W^* = \begin{pmatrix} w^* & cw^* & \cdots & c^{d-1}w^* \\ c^{-1}w^* & w^* & \cdots & c^{d-2}w^* \\ \vdots & \vdots & \ddots & \vdots \\ c^{1-d}w^* & c^{2-d}w^* & \cdots & w^* \end{pmatrix}, \quad c = O\left(\frac{1}{d}\right).$$

We pre-train only on the identical diagonal observations w^* using $\Omega_{\mathrm{pre}}^{(d)}$, with initialization $\boldsymbol{A}(0) = \boldsymbol{B}(0) = \alpha \boldsymbol{I}_d$ up to time T_1 (see Proposition 4.1 for the pre-training solution). We then reveal the full upper-triangular set $\Omega_{\mathrm{post}}^{(d)} = \{(i,j): 1 \leq i \leq j \leq d\}$ to form connectivity and continue training. By Theorem 4.3, for every $t \geq T_1$, the stable rank of $\boldsymbol{A}(t)$ is uniformly lower-bounded by $\Omega(d)$:

$$\frac{\|\boldsymbol{A}(t)\|_F^2}{\|\boldsymbol{A}(t)\|_2^2} \ge \left(\frac{4d-1}{4\sqrt{d}+1}\right)^2.$$

5 CONCLUSION

We demonstrate that in matrix completion, deeper networks ($L \geq 3$) inherently exhibit a stronger low-rank bias than shallow networks, primarily due to their coupled training dynamics, which operate regardless of observation patterns. For tractable analysis, we consider gradient flow starting at a family of deterministic initializations, showing in the diagonal observation setting that depth amplifies the low-rank bias. Furthermore, our theoretical analysis of warm-starting scenarios details the loss of plasticity phenomenon, revealing how large-norm, high-rank initial states can hinder convergence to low-rank solutions. We believe the theoretical results from matrix completion provide broader insight into how depth shapes implicit bias and explains the loss of plasticity in practical deep networks.

ETHICS STATEMENT

This work is purely theoretical and involves no human subjects, personal data, or new dataset collection. We foresee no safety, fairness, or privacy risks and confirm that we are in accordance with the ICLR Code of Ethics.

REPRODUCIBILITY STATEMENT

The proofs of all theorems and propositions in the main text appear in the corresponding appendices: Theorem 3.1 in Appendix D.1, Proposition 3.2 in Appendix D.2, Theorem 3.3 in Appendix D.3, Proposition 4.1 in Appendix E.1, and Theorems 4.2 and 4.3 in Appendices E.2 and E.3, respectively.

REFERENCES

- Alessandro Achille, Matteo Rovere, and Stefano Soatto. Critical learning periods in deep networks. In *International Conference on Learning Representations*, 2018.
- Maksym Andriushchenko, Dara Bahri, Hossein Mobahi, and Nicolas Flammarion. Sharpness-aware minimization leads to low-rank features. *Advances in Neural Information Processing Systems*, 36: 47032–47051, 2023.
- Sanjeev Arora, Nadav Cohen, and Elad Hazan. On the optimization of deep networks: Implicit acceleration by overparameterization. In Jennifer Dy and Andreas Krause (eds.), *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pp. 244–253. PMLR, 10–15 Jul 2018. URL https://proceedings.mlr.press/v80/arora18a.html.
- Sanjeev Arora, Nadav Cohen, Wei Hu, and Yuping Luo. Implicit regularization in deep matrix factorization. *Advances in Neural Information Processing Systems*, 32, 2019.
- Jordan Ash and Ryan P Adams. On warm-starting neural network training. *Advances in neural information processing systems*, 33:3884–3894, 2020.
- Zhiwei Bai, Jiajie Zhao, and Yaoyu Zhang. Connectivity shapes implicit regularization in matrix factorization models for matrix completion. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.
- Tudor Berariu, Wojciech Czarnecki, Soham De, Jorg Bornschein, Samuel Smith, Razvan Pascanu, and Claudia Clopath. A study on the plasticity of neural networks. *arXiv preprint arXiv:2106.00042*, 2021.
- E.J. Candes and T. Tao. Decoding by linear programming. *IEEE Transactions on Information Theory*, 51(12):4203–4215, 2005.
- Yihong Chen, Kelly Marchisio, Roberta Raileanu, David Adelani, Pontus Lars Erik Saito Stenetorp, Sebastian Riedel, and Mikel Artetxe. Improving language plasticity via pretraining with active forgetting. *Advances in Neural Information Processing Systems*, 36:31543–31557, 2023.
- Lenaic Chizat, Edouard Oyallon, and Francis Bach. On lazy training in differentiable programming. *Advances in neural information processing systems*, 32, 2019.
- Shibhansh Dohare, Richard S Sutton, and A Rupam Mahmood. Continual backprop: Stochastic gradient descent with persistent randomness. *arXiv* preprint arXiv:2108.06325, 2021.
- Spencer Frei, Gal Vardi, Peter Bartlett, Nathan Srebro, and Wei Hu. Implicit bias in leaky reLU networks trained on high-dimensional data. In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum?id=JpbLyEI5EwW.
- Tomer Galanti, Zachary S Siegel, Aparna Gupte, and Tomaso Poggio. SGD and weight decay provably induce a low-rank bias in neural networks, 2023. URL https://openreview.net/forum?id=N7Tv4aZ4Cyx.

Tomer Galanti, Zachary S Siegel, Aparna Gupte, and Tomaso A Poggio. SGD and weight decay secretly minimize the rank of your neural network. In *NeurIPS 2024 Workshop on Mathematics of Modern Machine Learning*, 2024. URL https://openreview.net/forum?id=xhW2WyPhRP.

- Daniel Gissin, Shai Shalev-Shwartz, and Amit Daniely. The implicit bias of depth: How incremental learning drives generalization. In *International Conference on Learning Representations*, 2020. URL https://openreview.net/forum?id=H1ljOnNFwB.
- Suriya Gunasekar, Blake E Woodworth, Srinadh Bhojanapalli, Behnam Neyshabur, and Nati Srebro. Implicit regularization in matrix factorization. *Advances in neural information processing systems*, 30, 2017.
- J Fernando Hernandez-Garcia, Shibhansh Dohare, Jun Luo, and Rich S Sutton. Reinitializing weights vs units for maintaining plasticity in neural networks. *arXiv preprint arXiv:2508.00212*, 2025.
- Minyoung Huh, Hossein Mobahi, Richard Zhang, Brian Cheung, Pulkit Agrawal, and Phillip Isola. The low-rank simplicity bias in deep networks. *arXiv preprint arXiv:2103.10427*, 2021.
- Xiangyun Hui, Xiaoxuan Ma, Yixuan Yang, and Song Li. The implicit regularization of gradient flow on separable datasets in relu networks. *Neurocomputing*, pp. 131367, 2025. ISSN 0925-2312. doi: https://doi.org/10.1016/j.neucom.2025.131367. URL https://www.sciencedirect.com/science/article/pii/S0925231225020399.
- Maximilian Igl, Gregory Farquhar, Jelena Luketina, Wendelin Boehmer, and Shimon Whiteson. Transient non-stationarity and generalisation in deep reinforcement learning. *arXiv* preprint *arXiv*:2006.05826, 2020.
- Arthur Jacot. Implicit bias of large depth networks: a notion of rank for nonlinear functions. *arXiv* preprint arXiv:2209.15055, 2022.
- Ziwei Ji and Matus Telgarsky. Gradient descent aligns the layers of deep linear networks. In *International Conference on Learning Representations*, 2019a. URL https://openreview.net/forum?id=HJflg30qKX.
- Ziwei Ji and Matus Telgarsky. The implicit bias of gradient descent on nonseparable data. In *Conference on learning theory*, pp. 1772–1798. PMLR, 2019b.
- Jikai Jin, Zhiyuan Li, Kaifeng Lyu, Simon Shaolei Du, and Jason D Lee. Understanding incremental learning of gradient descent: A fine-grained analysis of matrix sensing. In *International Conference* on Machine Learning, pp. 15200–15238. PMLR, 2023.
- Hyunji Jung, Hanseul Cho, and Chulhee Yun. Convergence and implicit bias of gradient descent on continual linear classification. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=DTqx3iqjkz.
- Daesung Kim and Hye Won Chung. Rank-1 matrix completion with gradient descent and small random initialization. *Advances in Neural Information Processing Systems*, 36:10530–10566, 2023.
- Jiyeon Kim, Hyunji Lee, Hyowon Cho, Joel Jang, Hyeonbin Hwang, Seungpil Won, Youbin Ahn, Dohaeng Lee, and Minjoon Seo. Knowledge entropy decay during language model pretraining hinders new knowledge acquisition. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=eHehzSDUFp.
- Michael Kleinman, Alessandro Achille, and Stefano Soatto. Critical learning periods emerge even in deep linear networks. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=Aq35gl2c1k.
- Yiwen Kou, Zixiang Chen, and Quanquan Gu. Implicit bias of gradient descent for two-layer relu and leaky relu networks on nearly-orthogonal data. *Advances in Neural Information Processing Systems*, 36:30167–30221, 2023.

- Saurabh Kumar, Henrik Marklund, and Benjamin Van Roy. Maintaining plasticity in continual learning via regenerative regularization. *arXiv preprint arXiv:2308.11958*, 2023.
 - Hojoon Lee, Hyeonseo Cho, Hyunseung Kim, Donghu Kim, Dugki Min, Jaegul Choo, and Clare Lyle. Slow and steady wins the race: Maintaining plasticity with hare and tortoise networks. In *ICML*, 2024. URL https://openreview.net/forum?id=VF177x7Syw.
 - Hojoon Lee, Dongyoon Hwang, Donghu Kim, Hyunseung Kim, Jun Jet Tai, Kaushik Subramanian, Peter R. Wurman, Jaegul Choo, Peter Stone, and Takuma Seno. Simba: Simplicity bias for scaling up parameters in deep reinforcement learning. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=jXLiDKsuDo.
 - Yuanzhi Li, Tengyu Ma, and Hongyang Zhang. Algorithmic regularization in over-parameterized matrix sensing and neural networks with quadratic activations. In *Conference On Learning Theory*, pp. 2–47. PMLR, 2018.
 - Zhiyuan Li, Yuping Luo, and Kaifeng Lyu. Towards resolving the implicit bias of gradient descent for matrix factorization: Greedy low-rank learning. In *International Conference on Learning Representations*, 2021. URL https://openreview.net/forum?id=AHOs7Sm5H7R.
 - Clare Lyle, Zeyu Zheng, Evgenii Nikishin, Bernardo Avila Pires, Razvan Pascanu, and Will Dabney. Understanding plasticity in neural networks. In *International Conference on Machine Learning*, pp. 23190–23211. PMLR, 2023.
 - Clare Lyle, Gharda Sokar, Razvan Pascanu, and Andras Gyorgy. What can grokking teach us about learning under nonstationarity? *arXiv preprint arXiv:2507.20057*, 2025.
 - Jianhao Ma and Salar Fattahi. Global convergence of sub-gradient method for robust matrix recovery: Small initialization, noisy measurements, and over-parameterization. *Journal of Machine Learning Research*, 24(96):1–84, 2023.
 - Jianhao Ma and Salar Fattahi. Convergence of gradient descent with small initialization for unregularized matrix completion. In *The Thirty Seventh Annual Conference on Learning Theory*, pp. 3683–3742. PMLR, 2024.
 - Govind Menon. The geometry of the deep linear network. arXiv preprint arXiv:2411.09004, 2024.
 - Behnam Neyshabur, Ryota Tomioka, and Nathan Srebro. In search of the real inductive bias: On the role of implicit regularization in deep learning. *arXiv preprint arXiv:1412.6614*, 2014.
 - Behnam Neyshabur, Ryota Tomioka, Ruslan Salakhutdinov, and Nathan Srebro. Geometry of optimization and implicit regularization in deep learning. *CoRR*, abs/1705.03071, 2017. URL http://arxiv.org/abs/1705.03071.
 - Evgenii Nikishin, Max Schwarzer, Pierluca D'Oro, Pierre-Luc Bacon, and Aaron Courville. The primacy bias in deep reinforcement learning. In *International conference on machine learning*, pp. 16828–16847. PMLR, 2022.
 - Sangyeon Park, Isaac Han, Seungwon Oh, and Kyung-Joong Kim. Activation by intervalwise dropout: A simple way to prevent neural networks from plasticity loss. *arXiv* preprint *arXiv*:2502.01342, 2025.
 - Noam Razin and Nadav Cohen. Implicit regularization in deep learning may not be explainable by norms. *Advances in neural information processing systems*, 33:21174–21187, 2020.
- Noam Razin, Asaf Maman, and Nadav Cohen. Implicit regularization in tensor factorization. In *International Conference on Machine Learning*, pp. 8913–8924. PMLR, 2021.
- Seyed Roozbeh Razavi Rohani, Khashayar Khajavi, Wesley Chung, Mo Chen, and Sharan Vaswani. Preserving plasticity in continual learning with adaptive linearity injection. *arXiv preprint arXiv:2505.09486*, 2025.
- Olivier Roy and Martin Vetterli. The effective rank: A measure of effective dimensionality. In 2007 15th European signal processing conference, pp. 606–610. IEEE, 2007.

- Baekrok Shin, Junsoo Oh, Hanseul Cho, and Chulhee Yun. Dash: Warm-starting neural network training in stationary settings without loss of plasticity. *Advances in Neural Information Processing Systems*, 37:43300–43340, 2024.
 - Mahdi Soltanolkotabi, Dominik Stöger, and Changzhi Xie. Implicit balancing and regularization: Generalization and convergence guarantees for overparameterized asymmetric matrix sensing. In *The Thirty Sixth Annual Conference on Learning Theory*, pp. 5140–5142. PMLR, 2023.
 - Daniel Soudry, Elad Hoffer, Mor Shpigel Nacson, Suriya Gunasekar, and Nathan Srebro. The implicit bias of gradient descent on separable data. *Journal of Machine Learning Research*, 19(70):1–57, 2018.
 - Jacob Mitchell Springer, Sachin Goyal, Kaiyue Wen, Tanishq Kumar, Xiang Yue, Sadhika Malladi, Graham Neubig, and Aditi Raghunathan. Overtrained language models are harder to fine-tune. In *Forty-second International Conference on Machine Learning*, 2025. URL https://openreview.net/forum?id=YW6edSufht.
 - Dominik Stöger and Mahdi Soltanolkotabi. Small random initialization is akin to spectral learning: Optimization and generalization guarantees for overparameterized low-rank matrix reconstruction. *Advances in Neural Information Processing Systems*, 34:23831–23843, 2021.
 - Matus Telgarsky. Deep learning theory lecture notes. *Lecture Notes v0. 0-e7150f2d (alpha), Univ. Illinois Urbana-Champaign, Champaign, IL, USA*, 2021.
 - Nadav Timor, Gal Vardi, and Ohad Shamir. Implicit regularization towards rank minimization in relu networks. In *International Conference on Algorithmic Learning Theory*, pp. 1429–1459. PMLR, 2023.
 - Simon Vock and Christian Meisel. Critical dynamics governs deep learning. *arXiv preprint* arXiv:2507.08527, 2025.
 - Blake Woodworth, Suriya Gunasekar, Jason D Lee, Edward Moroshko, Pedro Savarese, Itay Golan, Daniel Soudry, and Nathan Srebro. Kernel and rich regimes in overparametrized models. In *Conference on Learning Theory*, pp. 3635–3673. PMLR, 2020.
 - Chulhee Yun, Shankar Krishnan, and Hossein Mobahi. A unifying view on implicit bias in training linear neural networks. In *International Conference on Learning Representations*, 2021. URL https://openreview.net/forum?id=ZsZM-4iMQkH.
 - Chenyang Zhang, Difan Zou, and Yuan Cao. The implicit bias of adam on separable data. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL https://openreview.net/forum?id=xRQxan3WkM.
 - Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning requires rethinking generalization. In *International Conference on Learning Representations*, 2017. URL https://openreview.net/forum?id=Sy8gdB9xx.

CONTENTS Introduction **Problem Setting Implicit Bias of Depth Induced By Coupled Training Dynamics** Warm-up: Coupled Dynamics vs. Decoupled Dynamics in Depth-2 Networks . . . 3.2 Coupled Dynamics in Deep Networks Induce Implicit Bias Towards Low Rank . . . **Understanding Loss of Plasticity in Depth-2 Matrix Completion** 4.3 **Conclusion Further Related Works B** Coupled and Decoupled Training Dynamics C Additional Experiments D Proof for Section 3 E Proof for Section 4 **Useful Lemmas**

DECLARATION OF LLM USAGE

Large Language Models (LLM) were used solely to aid or polish writing. They did not generate ideas, analyses, or conclusions. All LLM-assisted text was reviewed and edited by the authors.

A FURTHER RELATED WORKS

A.1 IMPLICIT REGULARIZATION IN NEURAL NETWORKS

A substantial body of work investigates the *implicit regularization* of gradient-based training in overparameterized models; see, e.g., Gunasekar et al. (2017); Soudry et al. (2018); Woodworth et al. (2020); Yun et al. (2021); Ji & Telgarsky (2019a;b); Stöger & Soltanolkotabi (2021); Arora et al. (2019); Li et al. (2021); Andriushchenko et al. (2023); Jacot (2022); Frei et al. (2023); Kou et al. (2023); Timor et al. (2023); Zhang et al. (2024); Jung et al. (2025); Razin et al. (2021); Hui et al. (2025) for representative results. Within this line, we narrow our scope to the direction most relevant to our setting—how depth induces a bias toward low-rank solutions.

Several works investigate how depth promotes low-rank solutions (Gissin et al., 2020; Huh et al., 2021; Timor et al., 2023; Arora et al., 2019; Li et al., 2021). Huh et al. (2021) provide empirical evidence that deeper networks (both linear and nonlinear) tend to find solutions with lower effective-rank embeddings. Complementing this, Timor et al. (2023) show theoretically that ReLU networks trained with squared loss exhibit a bias toward low-rank solutions under the assumption that gradient flow converges to the solution minimizing the ℓ_2 norm.

Turning to deep linear networks, Gissin et al. (2020) and Li et al. (2021) study depth-induced bias as a function of initialization scale. They report that, as depth increases, the dependence on initialization can become weaker, and incremental learning can emerge. However, their analyses consider a matrix factorization task, which they frame as matrix completion with full observations. Therefore, in their setting, convergence to a low-rank solution is guaranteed if the model converges to zero-loss, which does not hold in our matrix completion task settings.

While Arora et al. (2019) investigate the matrix completion task in deep linear networks, offering insights from derived singular value dynamics, they cannot fully track these dynamics to prove low-rank convergence as network depth increases. Their analysis is primarily restricted to the regime where $t \geq t_0$, after which singular vectors are assumed to have stabilized. For $t \geq t_0$, they find that one singular value can be expressed as a function of another, involving a constant term that emerges from the state at t_0 (which can be the dominant component). Based on this derivation, they demonstrate that the gap between these singular values widens with increasing depth. In contrast, our Theorem 3.3, by precisely tracking the converged values of singular values, rigorously establishes their ultimate behavior and the resulting low-rank bias.

For depth-2 matrix completion tasks, Bai et al. (2024) introduce the connectivity argument. They prove that if the observations construct a connected bipartite graph, the model can converge to a low-rank solution when the initialization scale is infinitesimally small, subject to certain technical assumptions. Conversely, if the observations form a disconnected graph, the model generally cannot converge to a low-rank solution. However, a special case occurs if this disconnected graph is composed of complete bipartite components: here, the model converges to the minimum nuclear norm solution, again under specific technical assumptions. This characterization of implicit bias does not readily generalize to matrices with deeper matrices, as depicted in Figure 1.

A.2 Loss of Plasticity

Loss of plasticity describes a widely observed phenomenon where a model's ability to adapt to new information diminishes over time (Shin et al., 2024; Ash & Adams, 2020; Nikishin et al., 2022; Dohare et al., 2021; Achille et al., 2018; Lee et al., 2025; 2024; Lyle et al., 2025; Springer et al., 2025; Kim et al., 2025). The phenomenon is frequently observed in scenarios with gradually changing datasets, such as those encountered in reinforcement learning (Lyle et al., 2023; Nikishin et al., 2022; Igl et al., 2020) or continual learning (Kumar et al., 2023; Chen et al., 2023; Dohare et al., 2021; Park et al., 2025; Hernandez-Garcia et al., 2025; Rohani et al., 2025), where the model may struggle to adapt to new environments.

Although loss of plasticity is typically studied in non-stationary settings, a similar effect arises in stationary regimes where the dataset grows incrementally while the underlying distribution remains fixed (Shin et al., 2024; Ash & Adams, 2020; Berariu et al., 2021). In such cases, a model is first trained to convergence on an initial i.i.d. subset (e.g., a subset of CIFAR-10/100) and then warm-started for continued training on an expanded sample from the same distribution (e.g., the

full CIFAR-10/100). Perhaps counterintuitively, these warm-started models often generalize worse, yielding lower test accuracy than models trained from scratch on the combined dataset.

While this phenomenon is problematic in many real-world applications where new data is continuously added, theoretical studies on it remain scarce. Shin et al. (2024), for instance, offer a theoretical explanation using an artificial framework. Within this framework, they demonstrate that such behavior occurs because warm-started models often complete training by memorizing data-dependent noise, which is not useful for generalization. However, the analytical framework they employ is considered artificial and limited in its ability to accurately characterize the optimization processes of typical deep learning models.

Recently, Kleinman et al. (2024) observed loss of plasticity in deep linear networks, identifying "critical learning periods": an initial phase of effective learning followed by a significantly reduced capacity to learn later (Achille et al., 2018; Vock & Meisel, 2025). They employ a matrix completion framework to further observe this behavior. When observed that a model initially trained on a sparse set of observations and subsequently retrained (i.e., warm-started) on an expanded dataset typically exhibits a larger performance gap (in terms of reconstruction error) compared to a model trained from scratch on the entire expanded dataset. However, their work does not offer theoretical guarantees to account for these observations. Motivated by this, in Section 4, we attempt to explain this behavior within the specific context of depth-2 matrix completion settings.

B COUPLED AND DECOUPLED TRAINING DYNAMICS

This section introduces coupled and decoupled training dynamics (Definition 2) and illustrates them with concrete examples. Before that, we present Proposition B.1, which shows that for deep models $(L \ge 3)$, generic (absolutely continuous) initialization yields coupled dynamics almost surely.

Lemma B.1. Define $W_{b:a} \triangleq W_b W_{b-1} \cdots W_a$, and $W_{a:b} \triangleq I_d$ where $b \geq a$. For $w_{ij}(t) \triangleq e_i^\top W_{L:1}(t)e_j$,

$$\nabla_{\boldsymbol{W}_{l}} w_{ij}(t) = \left(\boldsymbol{W}_{L:l+1}(t)^{\top} \boldsymbol{e}_{i}\right) \left(\boldsymbol{W}_{l-1:1}(t) \boldsymbol{e}_{j}\right)^{\top} \in \mathbb{R}^{d \times d}.$$

Hence, for any (i, j) and (p, q),

$$\langle \nabla_{\boldsymbol{\theta}} w_{ij}(t), \nabla_{\boldsymbol{\theta}} w_{pq}(t) \rangle = \sum_{l=1}^{L} \left(\boldsymbol{e}_{i}^{\top} \boldsymbol{T}_{l}(t) \boldsymbol{e}_{p} \right) \left(\boldsymbol{e}_{j}^{\top} \boldsymbol{S}_{l}(t) \boldsymbol{e}_{q} \right),$$

where $T_l(t) \triangleq W_{L:l+1}(t)W_{L:l+1}(t)^{\top}$ and $S_l(t) \triangleq W_{l-1:1}(t)^{\top}W_{l-1:1}(t)$ are symmetric positive semidefinite matrix.

Proof. Define
$$a_l^{(i)}(t) \triangleq W_{L:l+1}(t)^{\top} e_i$$
 and $b_l^{(j)}(t) \triangleq W_{l-1:1}(t) e_j$. By

$$w_{ij}(t) = e_i^{\top} W_{L:l+1}(t) W_l(t) W_{l-1:1}(t) e_j = a_l^{(i)}(t)^{\top} W_l(t) b_l^{(j)}(t),$$

we have $\nabla_{\mathbf{W}_l} w_{ij}(t) = \mathbf{a}_l^{(i)}(t) \mathbf{b}_l^{(j)}(t)^{\top}$. Furthermore,

$$\begin{split} \langle \nabla_{\boldsymbol{\theta}} w_{ij}(t), \nabla_{\boldsymbol{\theta}} w_{pq}(t) \rangle &= \sum_{l=1}^{L} \left\langle \nabla_{\boldsymbol{W}_{l}} w_{ij}(t), \nabla_{\boldsymbol{W}_{l}} w_{pq}(t) \right\rangle_{F} \\ &= \sum_{l=1}^{L} \left\langle \boldsymbol{a}_{l}^{(i)}(t) \boldsymbol{b}_{l}^{(j)}(t)^{\top}, \boldsymbol{a}_{l}^{(p)}(t) \boldsymbol{b}_{l}^{(q)}(t)^{\top} \right\rangle_{F} \\ &= \sum_{i=1}^{L} \left(\boldsymbol{a}_{l}^{(i)}(t)^{\top} \boldsymbol{a}_{l}^{(p)}(t) \right) \left(\boldsymbol{b}_{l}^{(j)}(t)^{\top} \boldsymbol{b}_{l}^{(q)}(t) \right) \\ &= \sum_{i=1}^{L} \left(\boldsymbol{e}_{i}^{\top} \boldsymbol{T}_{l}(t) \boldsymbol{e}_{p} \right) \left(\boldsymbol{e}_{j}^{\top} \boldsymbol{S}_{l}(t) \boldsymbol{e}_{q} \right), \end{split}$$

which concludes the proof.

Proposition B.1. Let $L \geq 3$ and initialize $\{W_l(0)\}_{l=1}^L$ with i.i.d. entries from any absolutely continuous distribution. For any observation set $\Omega \subseteq [d] \times [d]$ where $|\Omega| \geq 2$, with probability 1,

$$\langle \nabla_{\boldsymbol{\theta}} w_{ij}(0), \nabla_{\boldsymbol{\theta}} w_{pq}(0) \rangle \neq 0$$

holds for all distinct $(i,j), (p,q) \in \Omega$. Consequently, no nontrivial partition $\Omega = \bigcup_{k=1}^K \Omega_k$ with $K \geq 2$ can satisfy the decoupling condition (6) at t = 0. Hence, by Definition 2, the gradient flow dynamics are coupled with probability I irrespective of the observation pattern.

Proof. By Lemma B.1, at t = 0 we have

$$arphi_{ij,pq}(oldsymbol{W}_1,\ldots,oldsymbol{W}_L) riangleq \left\langle
abla_{oldsymbol{ heta}} w_{ij},\,
abla_{oldsymbol{ heta}} w_{pq}
ight
angle = \sum_{l=1}^L \left(oldsymbol{e}_i^ op oldsymbol{T}_l oldsymbol{e}_p
ight) \left(oldsymbol{e}_j^ op oldsymbol{S}_l oldsymbol{e}_q
ight),$$

which is a polynomial in the entries of $\{W_l\}_{l=1}^L$. For any $(i,j) \neq (p,q)$, we now show that $\varphi_{ij,pq}$ is not the zero polynomial.

If i = p, the l = L term reduces to $e_j^{\top} S_L e_q$. By choosing $W_{1:L}$ so that S_L has a nonzero (j,q) entry, this term evaluates to a nonzero value; hence $\varphi_{ij,pq}$ is not identically zero. By symmetry, the same argument applies when j = q.

If $i \neq p$ and $j \neq q$, consider l = 2. Setting all other layers to I_d , choose W_3 so that $(e_i^\top T_2 e_p) \neq 0$ and choose W_1 so that $(e_j^\top S_2 e_q) \neq 0$. Then $\varphi_{ij,pq} = (e_i^\top T_2 e_p)(e_j^\top S_2 e_q) \neq 0$. Consequently, in all cases $\varphi_{ij,pq}$ is not identically zero.

Since $\varphi_{ij,pq}$ is a nonzero polynomial in the entries of $\{\boldsymbol{W}_l\}_{l=1}^L$, its zero set $Z_{ij,pq} \triangleq \{(\boldsymbol{W}_1,\ldots,\boldsymbol{W}_L): \varphi_{ij,pq}(\boldsymbol{W}_1,\ldots,\boldsymbol{W}_L)=0\}$ is a proper algebraic set in \mathbb{R}^{Ld^2} and hence has Lebesgue measure zero.

Let the initialization distribution of $(W_1(0), \dots, W_L(0))$ be absolutely continuous with respect to Lebesgue measure. Then

$$\Pr[(\mathbf{W}_1(0), \dots, \mathbf{W}_L(0)) \in Z_{ij,pq}] = 0,$$

so for this fixed pair $(i, j) \neq (p, q)$ we have $\varphi_{ij,pq}(\mathbf{W}_1(0), \dots, \mathbf{W}_L(0)) \neq 0$ almost surely. There are only finitely many distinct pairs in Ω . A finite union of measure-zero sets still has measure zero; hence, with probability one,

$$\varphi_{ij,pq} \neq 0 \quad \text{for all distinct } (i,j), (p,q) \in \Omega.$$
 (11)

By Definition 2, a decomposition $\Omega = \bigcup_{k=1}^K \Omega_k$ $(K \ge 2)$ yields decoupled dynamics only if $\langle \nabla_{\boldsymbol{\theta}} w_{ij}(t), \nabla_{\boldsymbol{\theta}} w_{pq}(t) \rangle = 0$

for all $(i, j) \in \Omega_k$, $(p, q) \in \Omega_l$ with $k \neq l$ and for all $t \geq 0$.

However, this already fails at t=0, since every cross-pair inner product is nonzero by (11). Thus, no such partition exists. Consequently, for $L \geq 3$ and any observation set Ω , the gradient flow dynamics are *coupled almost surely* under any absolutely continuous initialization.

B.1 COUPLED DYNAMICS EXAMPLE

B.1.1 Depth-2 Model

For shallow (L=2) matrices, coupled dynamics typically correspond to connected observations under generic initialization, in accordance with Definitions 1 and 2 (the specific case of initialization, such as zero matrices, which leads to decoupled dynamics, will be further detailed in a later subsection). We illustrate this principle with an example where the observed entries form the first column of a 2×2 matrix.

Consider a 2×2 matrix, denoted M_C , which is to be completed using its first column as observations:

$$M_{\mathrm{C}} \triangleq \begin{bmatrix} w_{11}^* & ? \\ w_{21}^* & ? \end{bmatrix}$$
.

The corresponding observation pattern matrix $P_{\rm C}$ is:

$$\mathbf{P}_{\mathrm{C}} = \begin{bmatrix} 1 & 0 \\ 1 & 0 \end{bmatrix}$$
.

The associated adjacency matrix $A_{\rm C}$ for the bipartite graph is constructed as:

$$\mathcal{A}_{\mathrm{C}} = egin{bmatrix} \mathbf{0}_{2,2} & \mathbf{P}_{\mathrm{C}}^{\top} \\ \mathbf{P}_{\mathrm{C}} & \mathbf{0}_{2,2} \end{bmatrix} = egin{bmatrix} 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 \end{bmatrix},$$

which forms a connected graph as illustrated in Figure 1a. This setup leads to coupled training dynamics under non-zero initialization. The coupling arises because parameters used to construct w_{11} and w_{21} overlap. Specifically, elements from the first column of matrix B (i.e., b_{11}, b_{21}) are common to the computation of both w_{11} and w_{21} . This shared dependency links the dynamics. The below illustration highlights these shared (teal) and distinct (red/blue) parameters involved in forming the observed entries w_{11} and w_{21} :

$$\begin{bmatrix} \mathbf{w}_{11} & \mathbf{w}_{12} \\ \mathbf{w}_{21} & \mathbf{w}_{22} \end{bmatrix} = \begin{bmatrix} \mathbf{a}_{11} & \mathbf{a}_{12} \\ \mathbf{a}_{21} & \mathbf{a}_{22} \end{bmatrix} \begin{bmatrix} b_{11} & b_{12} \\ b_{21} & b_{22} \end{bmatrix}$$
$$\mathbf{w}_{11} = \mathbf{a}_{11}b_{11} + \mathbf{a}_{12}b_{21}$$
$$\mathbf{w}_{21} = \mathbf{a}_{21}b_{11} + \mathbf{a}_{22}b_{21}$$

The shared use of b_{11} and b_{21} in reconstructing both observed entries is what couples their learning dynamics.

B.1.2 DEPTH≥ 3 MODEL

 For deeper matrices ($L \ge 3$), training dynamics are typically coupled, irrespective of the observation pattern (See Proposition B.1). Consider, for instance, predicting entries from the disconnected matrix M_D where only diagonal elements are observed:

$$M_{\mathrm{D}} \triangleq \begin{bmatrix} w_{11}^* & ? \\ ? & w_{22}^* \end{bmatrix}.$$

Even with such observations, for $L \geq 3$, coupling arises because parameters in intermediate layers are involved in computing multiple observed entries. This is illustrated in the following depth-3 example $(W_{3:1} = W_1W_2W_3)$. Elements of the intermediate matrix W_2 (colored teal) contribute to both the computation of w_{11} and w_{22} :

$$\begin{bmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \end{bmatrix} = \begin{bmatrix} (w_1)_{11} & (w_1)_{12} \\ (w_1)_{21} & (w_1)_{22} \end{bmatrix} \begin{bmatrix} (w_2)_{11} & (w_2)_{12} \\ (w_2)_{21} & (w_2)_{22} \end{bmatrix} \begin{bmatrix} (w_3)_{11} & (w_3)_{12} \\ (w_3)_{21} & (w_3)_{22} \end{bmatrix}.$$

Specifically, the observed entries are formed as:

$$w_{11} = \left((w_1)_{11} (w_2)_{11} + (w_1)_{12} (w_2)_{21} \right) (w_3)_{11}$$

$$+ \left((w_1)_{11} (w_2)_{12} + (w_1)_{12} (w_2)_{22} \right) (w_3)_{21},$$

$$w_{22} = \left((w_1)_{21} (w_2)_{11} + (w_1)_{22} (w_2)_{21} \right) (w_3)_{12}$$

$$+ \left((w_1)_{21} (w_2)_{12} + (w_1)_{22} (w_2)_{22} \right) (w_3)_{22}.$$

The shared involvement of all elements from W_2 (the teal matrix) in forming both w_{11} and w_{22} leads to coupled dynamics, provided these elements are non-zero. (Conversely, if some elements were to become zero, this could potentially lead to decoupled dynamics, as illustrated in the subsequent subsection.)

B.2 DECOUPLED DYNAMICS EXAMPLE

B.2.1 DEPTH-2 MODEL

For depth-2 models, decoupled dynamics coincide with disconnected observation patterns. Indeed, by Lemma B.1,

$$egin{aligned} \left\langle
abla_{m{ heta}} w_{ij},
abla_{m{ heta}} w_{pq}
ight
angle &= \sum_{l=1}^2 \left(oldsymbol{e}_i^{ op} oldsymbol{T}_l oldsymbol{e}_p
ight) \left(oldsymbol{e}_j^{ op} oldsymbol{S}_l oldsymbol{e}_q
ight) \ &= \left(oldsymbol{e}_1^{ op} oldsymbol{W}_2^{ op} oldsymbol{e}_p
ight) \delta_{jq} + \delta_{ip} \left(oldsymbol{e}_j^{ op} oldsymbol{W}_1^{ op} oldsymbol{W}_1 oldsymbol{e}_q
ight), \end{aligned}$$

where $\delta_{ab}=1$ if a=b and 0 otherwise. Hence, if $i\neq p$ and $j\neq p$, the inner product is identically zero for all weights, which explains the decoupling for the depth-2 matrix when the observations are disconnected.

To illustrate the disconnected case, consider the 2×2 incomplete matrix example M_D , to be completed from diagonal-only observations.

$$M_{\mathrm{D}} \triangleq \begin{bmatrix} w_{11}^* & ? \\ ? & w_{22}^* \end{bmatrix}.$$

Then the observation matrix P_D can be constructed as:

$$m{P}_{\! ext{D}} = egin{bmatrix} 1 & 0 \ 0 & 1 \end{bmatrix},$$

and the adjacency matrix A_D can be constructed as:

$$\mathcal{A}_{\mathrm{D}} = \begin{bmatrix} \mathbf{0}_{2,2} & \boldsymbol{P}_{\mathrm{D}}^{\top} \\ \boldsymbol{P}_{\mathrm{D}} & \mathbf{0}_{2,2} \end{bmatrix} = \begin{bmatrix} 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix},$$

which forms the disconnected graph as illustrated in Figure 1a. This setup inherently leads to decoupled training dynamics. The decoupling can be visually understood by examining how distinct sets of elements in the factor matrices A and B contribute to the observed entries w_{11} and w_{22} . Specifically, as illustrated below, red-colored entries are exclusively involved in predicting w_{11} , while blue-colored entries are exclusively involved in predicting w_{22} . These two sets of entries are disjoint, confirming the decoupled nature of the dynamics:

$$\begin{bmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \end{bmatrix} = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} \begin{bmatrix} b_{11} & b_{12} \\ b_{21} & b_{22} \end{bmatrix},$$

$$w_{11} = a_{11}b_{11} + a_{12}b_{21},$$

$$w_{22} = a_{21}b_{12} + a_{22}b_{22}.$$

B.2.2 DEPTH> 3 MODEL

For deep $(L \ge 3)$ matrices, decoupled training dynamics are observed in at least two key scenarios. First, as detailed in Appendix D.2.3, an αI_d initialization combined with diagonal-only observations leads to decoupled dynamics for any depth-factorized matrix.

To illustrate this for a deeper case, we revisit the $M_{\rm D}$ observation pattern in a depth-3 context. Lemma D.1 in Appendix D.2.3 states that with such an initialization and observing only diagonal entries, all off-diagonal elements of the factor matrices $W_l(t)$ remain zero throughout training. Consequently, the factor matrices W_1, W_2, W_3 are diagonal. The product matrix $W_{L:1}(t)$ is thus formed as:

$$\begin{bmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \end{bmatrix} = \begin{bmatrix} (w_1)_{11} & 0 \\ 0 & (w_1)_{22} \end{bmatrix} \begin{bmatrix} (w_2)_{11} & 0 \\ 0 & (w_2)_{22} \end{bmatrix} \begin{bmatrix} (w_3)_{11} & 0 \\ 0 & (w_3)_{22} \end{bmatrix}.$$

The observed entries are therefore computed as products of the respective diagonal elements:

$$w_{11} = (w_1)_{11}(w_2)_{11}(w_3)_{11},$$

$$w_{22} = (w_1)_{22}(w_2)_{22}(w_3)_{22}.$$

Since w_{11} depends only on the set of parameters $\{(\boldsymbol{W}_k)_{11}\}_{k=1}^3$ and w_{22} depends only on the entirely disjoint set of parameters $\{(\boldsymbol{W}_k)_{22}\}_{k=1}^3$, their training dynamics are decoupled.

Second, the training dynamics are also decoupled when all factor matrices are initialized as $d \times d$ zero matrices, $\mathbf{0}_{d \times d}$. To see this, note that by the chain rule, we have

$$\frac{\partial w_{pq}(t)}{\partial (w_l(t))_{ij}} = (\boldsymbol{W}_L(t)\boldsymbol{W}_{L-1}(t)\cdots\boldsymbol{W}_{l+1}(t))_{pi} (\boldsymbol{W}_{l-1}(t)\boldsymbol{W}_{l-2}(t)\cdots\boldsymbol{W}_{1}(t))_{jq}, \qquad (12)$$

where we define the (i, j)-th entry of the factor matrix $W_l(t) \triangleq (w_l(t))_{ij}$. If at some time t all factor matrices satisfy $W_k(t) = \mathbf{0}$, then the right-hand side of (12) is the zero matrix, and thus

$$\frac{\partial w_{pq}(t)}{\partial (w_l(t))_{ij}} = 0$$
 for all p, q .

Therefore,

$$\frac{\partial \phi}{\partial (w_l(t))_{ij}} = \sum_{(p,q) \in \Omega} \left(w_{pq}(t) - w_{pq}^* \right) \frac{\partial w_{pq}(t)}{\partial (W_l(t))_{ij}} = 0,$$

which implies

$$(w_l(t))_{ij} = -\frac{\partial \phi}{\partial (w_l(t))_{ij}} = 0.$$

Since the initial condition is $(w_l(0))_{ij} = 0$, uniqueness of ODE solutions guarantees that $(w_l(t))_{ij} \equiv 0$ for all $t \geq 0$. As this holds for arbitrary l, i, j, we conclude that $W_l(t) \equiv 0$ for all l and all $t \geq 0$.

Finally, because $\nabla_{\theta(t)} w_{pq}(t) = \mathbf{0}$ for all p, q and $t \geq 0$, the inner product condition

$$\langle \nabla_{\boldsymbol{\theta}(t)} w_{ij}(t), \nabla_{\boldsymbol{\theta}(t)} w_{pq}(t) \rangle = 0$$

is satisfied for all $(i,j),(p,q)\in\Omega$ and for all $t\geq 0$. Hence, the dynamics are (trivially) decoupled.

C ADDITIONAL EXPERIMENTS

This section provides additional experiments omitted from the main text.

C.1 IMPLICIT BIAS EXPERIMENTS

In Figure 1, we present experiments with specific choices of $M_{\rm C}$ and $M_{\rm D}$, which are 2×2 rank-1 ground-truth matrices illustrating connected and disconnected examples, respectively. To generalize these observations, we extended our experiments to a 3×3 rank-1 ground truth matrix, considering all possible connected and disconnected observation patterns. After accounting for symmetries to eliminate duplicates, this results in a total of 23 unique observation patterns, which are categorized into 17 connected and 6 disconnected cases.

For each of these 23 observation patterns, the 3×3 rank-1 ground truth matrix was generated using constituent vectors whose entries were sampled from a standard normal distribution. Each factor matrix was then initialized by sampling its entries from a Gaussian distribution with a mean of zero and a standard deviation of α . We performed 10 independent trials for each pattern.

Figure 4 illustrates that, consistent with the findings in Figure 1, a significant discrepancy exists between the behavior of depth-2 matrices and that of deeper matrices. This discrepancy becomes notably more pronounced for the disconnected observation patterns.

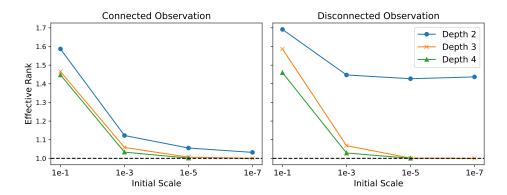


Figure 4: The left panel shows the averaged effective rank of all possible connected patterns as a function of the initial scale α^L . The right panel displays the averaged effective rank of all possible disconnected patterns.

We next provide a theoretical validation of our main claim: coupled dynamics induce a low-rank bias, whereas decoupled dynamics do not. This validation builds on Theorem 3.3, under various conditions, by numerically solving the equations while varying the ground truth value w^* and the dimension d. The results shown in Figure 7 (for $w^* = 1, d = 3$), Figure 5 (for $w^* = 10, d = 10$), and Figure 6 (for $w^* = 0.1, d = 10$) provide strong supporting evidence for the claim.

Furthermore, we ran gradient descent with a sufficiently small step size (10^{-5}) to validate our derived equations. For the results shown in Figure 8, we replicated the setup of Figure 7 ($w^*=1, d=3$), excluding the $\alpha=10^{-10}$ case due to prohibitive computation time. The observed values closely match the theoretical predictions from Theorem 3.3, as illustrated in Figure 7.

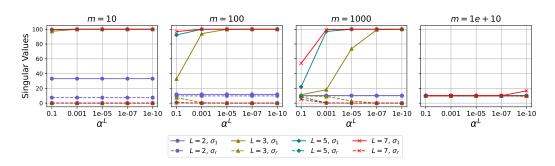


Figure 5: Numerical conditions identical to those in Figure 2, except with ground truth value $w^* = 10$ and dimension d = 10.

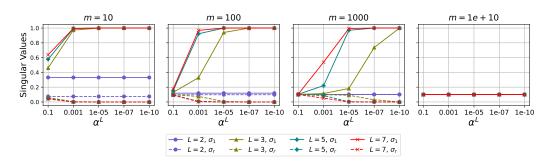


Figure 6: Numerical conditions identical to those in Figure 2, except with ground truth value $w^* = 0.1$ and dimension d = 10.

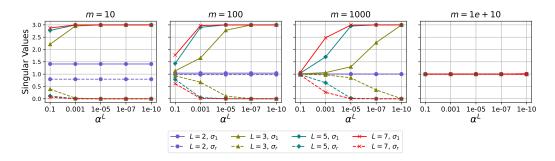


Figure 7: Numerical conditions identical to those in Figure 2, except with ground truth value $w^* = 1$ and dimension d = 3.

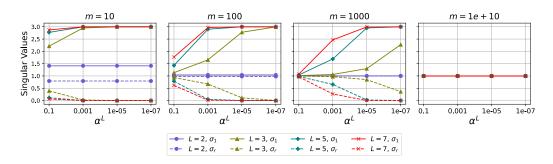
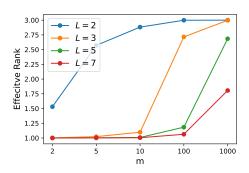
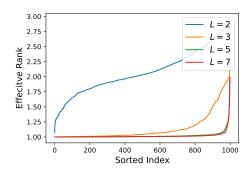


Figure 8: Gradient descent experiments conducted under conditions identical to those in Figure 7.





- (a) Results from Initialization using (7).
- (b) Results from Gaussian initialization.

Figure 9: (a) Effective rank for the initialization scheme in (7). The x-axis denotes the parameter m, which controls the initial rank characteristics of the model, while the y-axis represents the corresponding effective rank after convergence. (b) Effective rank distributions for Gaussian initialization. The results are from 1000 independent trials, sorted by their converged effective rank. The x-axis denotes the sorted trial index (from lowest to highest converged rank), and the y-axis represents the corresponding effective rank after convergence.

To validate that our initialization scheme (7) can achieve comparable outcomes to Gaussian initialization while offering more control, we conducted experiments on a 3×3 matrix completion task with diagonal observations (i.e., $w_{11}^*=w_{22}^*=w_{33}^*=1$). While our scheme allows initial rank properties to be adjusted via the parameter m, Gaussian initialization's inherent randomness precludes such direct control. Therefore, for comparison with Gaussian initialization, we ran 1000 trials (seeds) and sorted the converged solutions by their rank.

A comparison of the results in Figure 9 suggests that the behavioral trends may appear similar. In the depth-2 case, both initializations tend to converge to high-rank solutions. Moreover, for both initializations, a clear gap emerges between L=2 and L=3, with the depth-3 model exhibiting a stronger low-rank bias. For deeper networks ($L\geq 3$), the tendency to converge toward lower-rank solutions becomes increasingly pronounced as depth increases.

C.1.1 EXPERIMENTS IN NEURAL NETWORKS

To investigate the effect of depth on low-rank bias in practical settings, we train ResNet and VGG architectures across varying depths. While Huh et al. (2021) empirically show that deeper networks tend to produce embeddings of lower effective rank, their analysis focuses on feature embeddings rather than the weight matrices themselves. Therefore, following Galanti et al. (2023), we measure the effective rank of the weight matrices directly and find that deeper networks are biased toward low-rank solutions.

To be more specific, we train ResNet-18, 34, 50, and 101, as well as VGG-11, 13, 16, and 19, on CIFAR-10 and CIFAR-100 for 200 epochs with a batch size of 128. Training is performed using SGD with momentum 0.9, an initial learning rate of 0.1, weight decay of 0.0005, and a cosine annealing scheduler, together with standard data augmentation (horizontal flipping and random cropping). We measure the effective rank across all layers except the final one and average them to obtain a single scalar.

The results in Figures 10 to 13 show that the average effective rank decreases as depth increases. This observation aligns with Theorem 3.3, which provides a theoretical proof of the low-rank bias induced by depth in matrix completion settings.

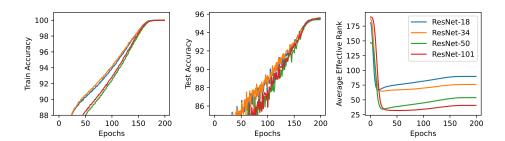


Figure 10: We train CIFAR-10 using ResNet models ranging from 18 to 101 layers, averaging results over five runs. The leftmost plot shows the training accuracy, the middle plot the test accuracy, and the rightmost plot the average effective rank. As depth increases, the average effective rank decreases.

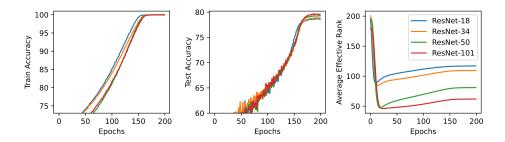


Figure 11: The results for CIFAR-100 with ResNet-18 to 101, under the same conditions as in Figure 10.

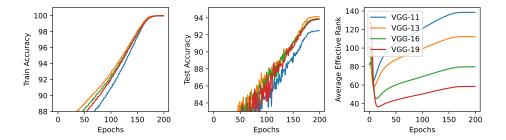


Figure 12: The results for CIFAR-10 with VGG-11 to 19, under the same conditions as in Figure 10.

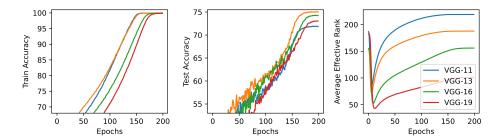


Figure 13: The results for CIFAR-100 with VGG-11 to 19, under the same conditions as in Figure 10.

C.2 Loss of Plasticity Experiments

Section 4.2 discusses a scenario where pre-training employs diagonal entries, after which an off-diagonal term (specifically, w_{12}^*) is introduced to restore connectivity, leading to coupled dynamics. Theorem 4.2 establishes that, in this situation, the model indeed does not converge to a low-rank solution. To empirically validate this theoretical finding, we conducted experiments using the family of initializations (7) tailored to this specific scenario, with results detailed in Figures 14 and 15. These experiments utilized a depth-2 model to reconstruct the ground-truth matrix, with an initialization scale set to $\alpha=10^{-35}$. Notably, if the initialization scale α is set significantly lower, as the dynamics are coupled, a cold-started model can converge to solutions exhibiting a more pronounced low-rank structure.

For the case presented in Figure 14, where $w^*=1, w_{12}^*=0.1$, following Theorem 4.2, the theoretical lower bound on the stable rank for a warm-started model initialized diagonally $(m=\infty)$ is approximately 1.45, while the empirically observed stable rank is approximately 1.8. Even in scenarios where substantial new information must be learned (e.g., by setting w_{12}^* to a large value), loss of plasticity is empirically observed, primarily manifesting as high test error (i.e., a significant gap between the target w_{21}^* and the converged w_{21}). While Theorem 4.2's analysis via stable rank does not fully explain an accompanying low-rank bias (a point consistent with Figure 15), the theorem does predict that w_{21} converges to a negative value, which implies a large test loss.

Furthermore, we performed additional experiments with different diagonal entry values to investigate whether this argument extends to other scenarios (results shown in Figure 16), although specific theoretical guarantees have not been established for these broader cases. We observe that even in these varied settings, both the effective rank and the stable rank of a warm-started model substantially exceed one, whereas cold-started models can converge to lower-rank solutions.

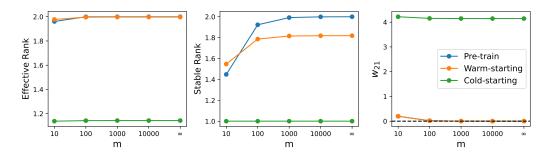


Figure 14: Experimental results for a 2×2 rank-1 ground-truth matrix W^* with $w_{11}^* = w_{22}^* = 1$ and $w_{12}^* = 0.5$ (implying $w_{21}^* = 2$ for rank-1 structure). Models, initialized according to (7), are first pre-trained on diagonal entries. After achieving zero-loss convergence in pre-training, the off-diagonal element w_{12}^* is introduced, and models are subsequently trained on combined diagonal and off-diagonal observations. The plots display: (Left and Middle) effective rank under different settings; (Right) converged value of $w_{21}(\infty)$. Key observations: (1) Warm-starting with a model that converged to a high-rank solution during pre-training tends to maintain this high rank, even when presented with the same subsequent observations as a cold-started model. (2) In the theoretically analyzed $m = \infty$ case, $w_{21}(\infty) < 0$ is observed, which correlates with the highest effective rank.

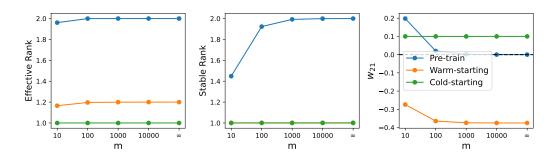


Figure 15: Experimental conditions identical to those in Figure 14, except with ground truth value $w_{12}^* = 10$. The model have to predict w_{21}^* as 0.1

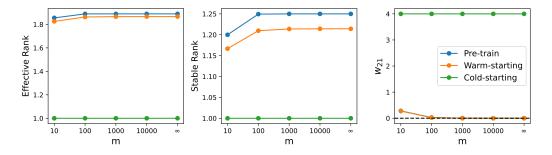


Figure 16: Experimental conditions identical to those in Figure 14, except with ground truth value $w_{11}^* = 1, w_{22}^* = 2$, and $w_{12}^* = 0.5$. The model have to predict w_{21}^* as 4.

D Proof for Section 3

 In this and the following sections, we prove the Propositions and Theorems presented in the main text. We begin with the proof of Theorem 3.1.

D.1 PROOF FOR THEOREM 3.1

When convergence is guaranteed, we can define the reference vector $\boldsymbol{u}^* \triangleq \frac{\boldsymbol{b}_1(\infty)}{\|\boldsymbol{b}_1(\infty)\|} \in \mathbb{R}^{d_1}$, which is entirely determined by their initial values and the targets. Note that \boldsymbol{u}^* does not change with time, since it is defined at $t = \infty$. We decompose $\boldsymbol{a}_1(t)$, $\boldsymbol{a}_2(t)$, and $\boldsymbol{b}_1(t)$ into two components: one parallel to \boldsymbol{u}^* and one perpendicular to \boldsymbol{u}^* :

$$a_1(t) = a_{1\parallel}(t) + a_{1\perp}(t), \quad a_2(t) = a_{2\parallel}(t) + a_{2\perp}(t), \quad b_1(t) = b_{1\parallel}(t) + b_{1\perp}(t).$$

For any vector $u \in \mathbb{R}^{d_1}$, the parallel component is defined as $u_{\parallel} = (u^*^{\top}u)u^*$, and the perpendicular component as $u_{\perp} = u - u_{\parallel}$.

We introduce notation to quantify the alignment of each vector with u^* :

$$\alpha_{\boldsymbol{a}_1}(t) = \boldsymbol{u}^{*\top} \boldsymbol{a}_1(t), \quad \alpha_{\boldsymbol{a}_2}(t) = \boldsymbol{u}^{*\top} \boldsymbol{a}_2(t), \quad \alpha_{\boldsymbol{b}_1}(t) = \boldsymbol{u}^{*\top} \boldsymbol{b}_1(t).$$
 (13)

Additionally, we define notation to measure the magnitude of the perpendicular components:

$$\beta_{\mathbf{a}_1}(t) = \|\mathbf{a}_{1\perp}(t)\|_2^2, \quad \beta_{\mathbf{a}_2}(t) = \|\mathbf{a}_{2\perp}(t)\|_2^2, \quad \beta_{\mathbf{b}_1}(t) = \|\mathbf{b}_{1\perp}(t)\|_2^2. \tag{14}$$

Then, using equation (4), time evolution of each component in equation (13) can be written as:

$$\alpha_{\boldsymbol{a}_{1}}(t) = \boldsymbol{u}^{*\top} \boldsymbol{a}_{1}(t)$$

$$= \underbrace{(\boldsymbol{w}_{11}^{*} - \boldsymbol{a}_{1}^{\top}(t)\boldsymbol{b}_{1}(t))}_{\triangleq r_{1}(t)} \boldsymbol{u}^{*\top} \boldsymbol{b}_{1}(t)$$

$$= r_{1}(t)\alpha_{\boldsymbol{b}_{1}}(t). \tag{15}$$

Likewise, for $\alpha_{a_2}(t)$, we derive:

$$\alpha_{\mathbf{a}_{2}}(t) = \mathbf{u}^{*\top} \mathbf{a}_{2}(t)
= \underbrace{(\mathbf{w}_{21}^{*} - \mathbf{a}_{2}^{\top}(t)\mathbf{b}_{1}(t))}_{\triangleq r_{2}(t)} \mathbf{u}^{*\top} \mathbf{b}_{1}(t)
= r_{2}(t)\alpha_{\mathbf{b}_{1}}(t).$$
(16)

Finally, for $\alpha_{b_1}(t)$, we have:

$$\alpha \dot{\boldsymbol{b}}_{1}(t) = \boldsymbol{u}^{*\top} \dot{\boldsymbol{b}}_{1}(t)
= (w_{11}^{*} - \boldsymbol{a}_{1}^{\top}(t)\boldsymbol{b}_{1}(t))\boldsymbol{u}^{*\top} \boldsymbol{a}_{1}(t) + (w_{21}^{*} - \boldsymbol{a}_{2}^{\top}(t)\boldsymbol{b}_{1}(t))\boldsymbol{u}^{*\top} \boldsymbol{a}_{2}(t)
= r_{1}(t)\alpha_{\boldsymbol{a}_{1}}(t) + r_{2}(t)\alpha_{\boldsymbol{a}_{2}}(t).$$
(17)

Also, for the perpendicular components, their time evolution can be derived as:

$$\begin{split} \dot{\beta}_{\boldsymbol{a}_{1}}(t) &= 2\boldsymbol{a}_{1\perp}(t) \cdot \dot{\boldsymbol{a}}_{1\perp}(t) \\ &= 2\boldsymbol{a}_{1\perp}(t) \cdot \frac{d}{dt} \left(\boldsymbol{a}_{1}(t) - \left(\boldsymbol{u}^{*\top} \boldsymbol{a}_{1}(t) \right) \boldsymbol{u}^{*} \right) \\ &= 2\boldsymbol{a}_{1\perp}(t) \cdot \left(r_{1}(t)\boldsymbol{b}_{1}(t) - r_{1}(t) \left(\boldsymbol{u}^{*\top} \boldsymbol{b}_{1}(t) \right) \boldsymbol{u}^{*} \right). \end{split}$$

Noting that $a_{1\perp}(t)$ is perpendicular to u^* , the second term in the parenthesis is zero. Thus, we have

$$\dot{\beta}_{\boldsymbol{a}_1}(t) = 2r_1(t)\boldsymbol{a}_{1\perp}(t)^{\top}\boldsymbol{b}_{1\perp}(t).$$

Likewise, for $\beta_{a_2}(t)$ and $\beta_{b_1}(t)$, we can derive their time derivative as:

$$\dot{\beta}_{a_2}(t) = 2r_2(t)a_{2\perp}(t)^{\top}b_{1\perp}(t), \quad \dot{\beta}_{b_1}(t) = \dot{\beta}_{a_1}(t) + \dot{\beta}_{a_2}(t).$$

Note that by the definition of u^* , we have $\beta_{b_1}(\infty) = 0$. Integrating the identity $\dot{\beta}_{b_1}(t) = \dot{\beta}_{a_1}(t) + \dot{\beta}_{a_2}(t)$ from t = 0 to ∞ gives:

$$\beta_{\boldsymbol{a}_1}(\infty) + \beta_{\boldsymbol{a}_2}(\infty) = \underbrace{\beta_{\boldsymbol{a}_1}(0) + \beta_{\boldsymbol{a}_2}(0) - \beta_{\boldsymbol{b}_1}(0)}_{\triangleq \beta_0 > 0}.$$

This equation shows that if the initial value β_0 is small, it constrains the total perpendicular magnitude at convergence. However, since we do not know u^* in advance, one natural way to ensure small perpendicular components is to initialize the entire norms of $a_1(0)$, $a_2(0)$ to be sufficiently small.

To develop a more rigorous understanding, we analyze the parallel components. Under the assumption of convergence, we have:

$$a_1(\infty)^{\top} b_1(\infty) = w_{11}^*, \quad a_2(\infty)^{\top} b_1(\infty) = w_{21}^*.$$

Decomposing $a_1(\infty)$ and $a_2(\infty)$ leads to:

$$\boldsymbol{a}_{1}(\infty)^{\top}\boldsymbol{b}_{1}(\infty) = \left(\boldsymbol{a}_{1\perp}(\infty) + \boldsymbol{u}^{*\top}\boldsymbol{a}_{1}(\infty)\boldsymbol{u}^{*}\right)^{\top}\boldsymbol{b}_{1}(\infty)$$
$$= \alpha_{\boldsymbol{a}_{1}}(\infty)\alpha_{\boldsymbol{b}_{1}}(\infty) = w_{11}^{*}, \tag{18}$$

$$\boldsymbol{a}_{2}(\infty)^{\top}\boldsymbol{b}_{1}(\infty) = \left(\boldsymbol{a}_{2\perp}(\infty) + \boldsymbol{u}^{*\top}\boldsymbol{a}_{2}(\infty)\boldsymbol{u}^{*}\right)^{\top}\boldsymbol{b}_{1}(\infty)$$
$$= \alpha_{\boldsymbol{a}_{2}}(\infty)\alpha_{\boldsymbol{b}_{1}}(\infty) = w_{21}^{*}. \tag{19}$$

Using equations (15)–(17), and noting that

$$\frac{d}{dt}\alpha_{\boldsymbol{b}_1}^2(t) = \frac{d}{dt}(\alpha_{\boldsymbol{a}_1}^2(t) + \alpha_{\boldsymbol{a}_2}^2(t)),$$

we can integrate both sides of the equation over time from 0 to ∞ to obtain:

$$\alpha_{a_1}^2(\infty) + \alpha_{a_2}^2(\infty) = \alpha_{b_1}^2(\infty) + \underbrace{\alpha_{a_1}^2(0) + \alpha_{a_2}^2(0) - \alpha_{b_1}^2(0)}_{\triangleq \alpha_0}.$$
 (20)

By solving equations (18), (19), and (20), we can obtain closed-form solutions of $\alpha_{a_1}(\infty)$, $\alpha_{a_2}(\infty)$, and $\alpha_{b_1}(\infty)$ as follows:

$$\alpha_{\mathbf{a}_{1}}^{2}(\infty) = \frac{2w_{11}^{*2}}{\sqrt{\alpha_{0}^{2} + 4w_{11}^{*2} + 4w_{21}^{*2} - \alpha_{0}}}, \quad \alpha_{\mathbf{a}_{2}}^{2}(\infty) = \frac{2w_{21}^{*2}}{\sqrt{\alpha_{0}^{2} + 4w_{11}^{*2} + 4w_{21}^{*2} - \alpha_{0}}}, \quad (21)$$

$$\alpha_{b_1}^2(\infty) = \frac{\sqrt{\alpha_0^2 + 4w_{11}^{*2} + 4w_{21}^{*2} - \alpha_0}}{2}.$$
 (22)

Thus, we can upper bound the proportion of the perpendicular component of $a_1(\infty)$ and $a_2(\infty)$ relative to its total magnitude as follows:

$$\frac{\|\boldsymbol{a}_{1\perp}(\infty)\|^2}{\|\boldsymbol{a}_{1}(\infty)\|^2} = \frac{\beta_{\boldsymbol{a}_{1}}(\infty)}{\alpha_{\boldsymbol{a}_{1}}^2(\infty) + \beta_{\boldsymbol{a}_{1}}(\infty)} \leq \frac{\beta_{0}\left(\sqrt{\alpha_{0}^2 + 4w_{11}^{*2} + 4w_{21}^{*2}} - \alpha_{0}\right)}{2w_{11}^{*2}},$$

$$\frac{\|\boldsymbol{a}_{2\perp}(\infty)\|^2}{\|\boldsymbol{a}_{2}(\infty)\|^2} = \frac{\beta_{\boldsymbol{a}_{2}}(\infty)}{\alpha_{\boldsymbol{a}_{2}}^2(\infty) + \beta_{\boldsymbol{a}_{2}}(\infty)} \leq \frac{\beta_{0}\left(\sqrt{\alpha_{0}^2 + 4w_{11}^{*2} + 4w_{21}^{*2}} - \alpha_{0}\right)}{2w_{21}^{*2}}.$$

To further refine these bounds, we analyze the terms β_0 and $S(\alpha_0) \triangleq \sqrt{\alpha_0^2 + 4w_{11}^*^2 + 4w_{21}^*^2} - \alpha_0$. By the definition of β_0 , it is upper bounded by $\|\boldsymbol{a}_1(0)\|^2 + \|\boldsymbol{a}_2(0)\|^2 = \|\boldsymbol{A}(0)\|_F^2$. Also, by the definition of α_0 , we have:

$$-\|\boldsymbol{b}_1(0)\|_2^2 \le \alpha_0 \le \|\boldsymbol{A}(0)\|_F^2.$$

Noting that the function $f(x) = \sqrt{x^2 + C} - x$ (where C > 0) is non-negative and monotonically decreasing for all $x \in \mathbb{R}$, we can upper bound $S(\alpha_0)$ using the lower bound of α_0 :

$$S(\alpha_0) \leq S(-\|\boldsymbol{b}_1(0)\|_2^2)$$

$$= \sqrt{(-\|\boldsymbol{b}_1(0)\|_2^2)^2 + 4(w_{11}^{*2} + w_{21}^{*2})} - (-\|\boldsymbol{b}_1(0)\|_2^2)$$

$$= \sqrt{\|\boldsymbol{b}_1(0)\|_2^4 + 4(w_{11}^{*2} + w_{21}^{*2})} + \|\boldsymbol{b}_1(0)\|_2^2.$$

Substituting these bounds for β_0 and $S(\alpha_0)$ into the inequality $\frac{\|\mathbf{a}_{1_{\perp}}(\infty)\|^2}{\|\mathbf{a}_{1}(\infty)\|_2^2} \leq \frac{\beta_0 S(\alpha_0)}{2w_{11}^*}$, we obtain the final upper bound for the proportion of the perpendicular component of $\mathbf{a}_1(\infty)$:

$$\frac{\|\boldsymbol{a}_{1\perp}(\infty)\|^2}{\|\boldsymbol{a}_{1}(\infty)\|_2^2} \leq \frac{\|\boldsymbol{A}(0)\|_F^2 \left(\sqrt{\|\boldsymbol{b}_{1}(0)\|_2^4 + 4(w_{11}^{*2} + w_{21}^{*2})} + \|\boldsymbol{b}_{1}(0)\|_2^2\right)}{2w_{11}^{*2}}.$$

A similar bound applies to $\frac{\|\boldsymbol{a}_{2\perp}(\infty)\|^2}{\|\boldsymbol{a}_{2}(\infty)\|_2^2}$:

$$\frac{\|\boldsymbol{a}_{2\perp}(\infty)\|^2}{\|\boldsymbol{a}_{2}(\infty)\|_2^2} \leq \frac{\|\boldsymbol{A}(0)\|_F^2 \left(\sqrt{\|\boldsymbol{b}_{1}(0)\|_2^4 + 4(w_{11}^{*~2} + w_{21}^{*~2})} + \|\boldsymbol{b}_{1}(0)\|_2^2\right)}{2w_{21}^{*~2}}.$$

D.2 Proof for Proposition 3.2

According to the definition of coupled/decoupled dynamics presented in Definition 2, for the family of initializations defined in (7) along with the diagonal observations ($\Omega_{\rm diag}^{(d)}$), we divide the cases to ensure that all possible scenarios for this family of initializations are covered.

D.2.1 Case for L=2

First, we consider the depth-2 (L=2) case. Each diagonal observation, $w_{ii}(t)$, is the inner product of the i-th row of $\boldsymbol{A}(t)$ and the i-th column of $\boldsymbol{B}(t)$. Then, when we take the gradient $\nabla_{\theta(t)}w_{ii}(t)$, where $\theta(t)$ represents the concatenation of $\boldsymbol{A}(t)$ and $\boldsymbol{B}(t)$, this gradient has non-zero components only corresponding to the i-th row of $\boldsymbol{A}(t)$ and the i-th column of $\boldsymbol{B}(t)$; all other components are zero for all $t \geq 0$. Therefore, for any $j \neq i$, the inner product $\langle \nabla_{\theta(t)}w_{ii}(t), \nabla_{\theta(t)}w_{jj}(t) \rangle$ must be zero. This means that there exists a partition of $\Omega_{\text{diag}}^{(d)}$ into disjoint subsets $\Omega_1, \ldots, \Omega_d$, where each $\Omega_i = \{(i,i)\}$. Therefore, for any initialization, the training dynamics are **decoupled**.

D.2.2 Case for $L \ge 3$ and $1 < m < \infty$

For the deeper matrix case $(L \geq 3)$, we first note that each diagonal observation $w_{ii}(t)$ can be expressed as:

$$w_{ii}(t) = \sum_{i_{L-1}=1}^{d} \cdots \sum_{i_{1}=1}^{d} (\boldsymbol{W}_{L}(t))_{i,i_{L-1}} (\boldsymbol{W}_{L-1}(t))_{i_{L-1},i_{L-2}} \cdots (\boldsymbol{W}_{1}(t))_{i_{1},i}.$$

Now consider the case $1 < m < \infty$, where every entry of each weight matrix $W_l(0)$ (for $l=1,\ldots,L$) is initialized as a positive value. Since $w_{ii}(0)$ is a sum of products of these positive entries, its gradient with respect to the parameters $\theta(0)$, $\nabla_{\theta(0)}w_{ii}(0)$, likewise consist of components that are sums of positive products (see (23)). Therefore, it is asserted that each relevant component of $\nabla_{\theta(0)}w_{ii}(0)$ is positive at initialization. Consequently, for any $j \neq i$, since both $\nabla_{\theta(0)}w_{ii}(0)$ and $\nabla_{\theta(0)}w_{jj}(0)$ have all their corresponding components positive, their inner product $\langle \nabla_{\theta(0)}w_{ii}(0), \nabla_{\theta(0)}w_{jj}(0) \rangle$ will be non-zero (specifically, positive). This non-zero inner product signifies **coupled dynamics**.

D.2.3 Case for $L \geq 3$ and $m = \infty$

Next, we examine the $m=\infty$ case, which corresponds to initializing each factor matrix $\mathbf{W}_l(0)$ as a scaled identity, i.e., $\mathbf{W}_l(0)=\alpha\mathbf{I}_d$. The following lemma states that under this initialization, and for dynamics driven by diagonal observations (from $\Omega_{\mathrm{diag}}^{(d)}$), all off-diagonal elements of each $\mathbf{W}_l(t)$ remain zero for all $t\geq 0$.

Lemma D.1. For a set of L matrices $W_1(t), \ldots, W_L(t) \in \mathbb{R}^{d \times d}$, let $W_{L:1}(t) = W_L(t) \cdots W_1(t)$. Following gradient flow dynamics in (3), if each factor matrix $W_l(0)$ is initialized as a diagonal matrix (e.g., $W_l(0) = \alpha_l I_d$ for scalars α_l), then all off-diagonal elements of each matrix $W_l(t)$ remain zero for all $t \geq 0$.

Proof. For a given diagonal observation indices $\Omega_{\mathrm{diag}}^{(d)}$, if we consider the gradient flow dynamics for an (i,j)-th entry of the factor matrix $W_l(t)$ ($\triangleq (w_l(t))_{ij}$), we have:

$$\frac{d(w_l(t))_{ij}}{dt} = -\frac{\partial \phi}{\partial (w_l(t))_{ij}}$$
$$= -\sum_{p=1}^d (w_{pp}(t) - w_{pp}^*) \frac{\partial w_{pp}(t)}{\partial (w_l(t))_{ij}},$$

Here, the derivative of a diagonal element $w_{pp}(t)$ with respect to $(w_l(t))_{ij}$ is:

$$\frac{\partial w_{pp}(t)}{\partial (w_l(t))_{ij}} = (\boldsymbol{W}_L(t)\boldsymbol{W}_{L-1}(t)\cdots\boldsymbol{W}_{l+1}(t))_{pi} (\boldsymbol{W}_{l-1}(t)\boldsymbol{W}_{l-2}(t)\cdots\boldsymbol{W}_{1}(t))_{jp}, \qquad (23)$$

where the first term is (p, i)-th element of the product $W_L(t)W_{L-1}(t)\cdots W_{l+1}(t)$, and the second term is (j, p)-th element of the product $W_{l-1}(t)W_{l-2}(t)\cdots W_1(t)$. We want to show that if all $W_l(t)$ are diagonal, then $\frac{d(w_l(t))_{ij}}{dt}=0$ for any off-diagonal element $(w_l(t))_{ij}$ (i.e., $i\neq j$).

Assume at a given time t that all factor matrices $W_l(t)$ are diagonal. Then, the product $P(t) \triangleq \prod_{k=l+1}^L W_k(t)$ is diagonal. Similarly, the product $S(t) \triangleq \prod_{k=1}^{l-1} W_k(t)$ is diagonal. For $\frac{\partial w_{pp}(t)}{\partial (w_l(t))_{ij}}$ to be non-zero (given all $W_l(t)$ are diagonal), both $(P(t))_{pi}$ and $(S(t))_{jp}$ must be non-zero. This requires p=i and j=p, which implies i=j.

However, we are considering an off-diagonal element $(w_l(t))_{ij}$, for which $i \neq j$. This means that if all $W_l(t)$ are diagonal, then for any p:

$$\frac{\partial w_{pp}}{\partial (w_l(t))_{ij}} = 0, \quad \text{if } i \neq j$$

Substituting this into the dynamic equation for $(w_l(t))_{ij}$:

$$\frac{d(w_l(t))_{ij}}{dt} = -\sum_{p=1}^{d} (w_{pp}(t) - w_{pp}^*) \cdot 0 = 0, \quad \text{if } i \neq j$$

Initially, $W_l(0)$ are diagonal, so all off-diagonal elements $(w_l(t))_{ij}$ are zero for $i \neq j$. Since their time derivatives are zero when they are zero (i.e., when the matrices are diagonal), these off-diagonal elements remain zero for all $t \geq 0$.

With Lemma D.1, the factor matrices $W_l(t)$ remain diagonal, so $w_{ii}(t) = (W_L(t))_{ii} \cdots (W_1(t))_{ii}$. This structure leads to decoupled dynamics because each $w_{ii}(t)$ depends exclusively on the set of parameters $\{(W_k(t))_{ii}\}_{k=1}^L$, while $w_{jj}(t)$ (for $j \neq i$) depends on the distinct set $\{(W_k(t))_{jj}\}_{k=1}^L$. Consequently, for any $j \neq i$, their respective gradients $\nabla_{\theta(t)}w_{ii}(t)$ and $\nabla_{\theta(t)}w_{jj}(t)$ are orthogonal, meaning their inner product is zero:

$$\langle \nabla_{\theta(t)} w_{ii}(t), \nabla_{\theta(t)} w_{jj}(t) \rangle = 0.$$

This orthogonality implies that the learning for each diagonal entry is independent, allowing a conceptual partition of $\Omega_{\rm diag}^{(d)}$ into disjoint subsets $\Omega_i = \{(i,i)\}$. Therefore, under this specific diagonal initialization (the $m=\infty$ case), the training dynamics are **decoupled**.

D.3 Proof for Theorem 3.3

Before presenting the proof of Theorem 3.3, we first restate the problem setting. The model is defined as $W_{L:1}(t) = W_L(t)W_{L-1}(t)\cdots W_1(t)$, where each factor matrix $W_l(t) \in \mathbb{R}^{d \times d}$ is subject to diagonal observations $\Omega_{\mathrm{diag}}^{(d)} = \{(i,i)\}_{i=1}^d$, and follows the gradient flow described in (3). We also assume that all diagonal entries are equal, i.e., $w^* \triangleq w_{11}^* = w_{22}^* \cdots = w_{dd}^*$. To simplify notation, we use $\ell(W_{L:1}(t))$ in place of $\ell(W_{L:1}(t); \Omega_{\mathrm{diag}}^{(d)})$ when the context is clear. The explicit gradient flow dynamics for each factor matrix is then given by:

$$\dot{\mathbf{W}}_{l}(t) = -\prod_{i=l+1}^{L} \mathbf{W}_{i}(t)^{\top} \cdot \nabla \ell(\mathbf{W}_{L:1}(t)) \cdot \prod_{i=1}^{l-1} \mathbf{W}_{i}(t)^{\top}, \tag{24}$$

where $\nabla \ell(\boldsymbol{W}_{L:1}(t)) = \operatorname{diag}(r_1(t), r_2(t), \cdots, r_d(t))$. Here, the residual term is defined as $r_i(t) \triangleq w_{ii}(t) - w^*$. To begin, we first present the preliminary lemma required for the following result.

Lemma D.2. Let I_n denote the $n \times n$ identity matrix and $J_n \triangleq \mathbb{1}_n \mathbb{1}_n^\top$ denote the $n \times n$ matrix with all entries equal to 1. Then the set

$$\mathcal{S} = \{ a \mathbf{I}_n + b \mathbf{J}_n \mid a, b \in \mathbb{R} \}$$

is closed under scalar multiplication, addition, and matrix multiplication. Also, any two matrices $A, B \in \mathcal{S}$ commute.

Proof. Let

$$\mathbf{A} = a\mathbf{I}_n + b\mathbf{J}_n$$
 and $\mathbf{B} = c\mathbf{I}_n + d\mathbf{J}_n$,

with $a, b, c, d \in \mathbb{R}$, and let $\lambda \in \mathbb{R}$ be an arbitrary scalar.

Scalar Multiplication:

$$\lambda \mathbf{A} = \lambda (a \mathbf{I}_n + b \mathbf{J}_n) = (\lambda a) \mathbf{I}_n + (\lambda b) \mathbf{J}_n.$$

Since $\lambda a, \lambda b \in \mathbb{R}$, it follows that $\lambda A \in \mathcal{S}$.

Addition:

$$\mathbf{A} + \mathbf{B} = (a\mathbf{I}_n + b\mathbf{J}_n) + (c\mathbf{I}_n + d\mathbf{J}_n) = (a+c)\mathbf{I}_n + (b+d)\mathbf{J}_n.$$

Since a + c, $b + d \in \mathbb{R}$, we have $A + B \in \mathcal{S}$.

Matrix Multiplication:

$$\mathbf{AB} = (a\mathbf{I}_n + b\mathbf{J}_n)(c\mathbf{I}_n + d\mathbf{J}_n).$$

Using the distributive property and the facts that

$$I_n J_n = J_n I_n = J_n$$
 and $J_n^2 = n J_n$,

we expand:

$$AB = ac \mathbf{I}_n \mathbf{I}_n + ad \mathbf{I}_n \mathbf{J}_n + bc \mathbf{J}_n \mathbf{I}_n + bd \mathbf{J}_n^2$$

= $ac \mathbf{I}_n + ad \mathbf{J}_n + bc \mathbf{J}_n + bd (n\mathbf{J}_n)$
= $ac \mathbf{I}_n + (ad + bc + nbd) \mathbf{J}_n$.

Thus, AB is of the form $\alpha I_n + \beta J_n$ with $\alpha = ac$ and $\beta = ad + bc + nbd$, and hence $AB \in \mathcal{S}$.

Commutativity: By the same procedure as above,

$$AB = (aI_n + bJ_n)(cI_n + dJ_n)$$

$$= acI_n + (ad + bc + nbd)J_n$$

$$= caI_n + (cb + da + ndb)J_n$$

$$= BA.$$

which completes the proof.

D.3.1 Case for L=2 & $L\geq 3$ and $1< m<\infty$

 We will first examine two main scenarios: the depth-2 (L=2) case and deeper networks $(L \ge 3)$ where $1 < m < \infty$. The $m = \infty$ case will be considered separately in the later subsection, as its initialization with αI_d warrants distinct treatment.

We now proceed to prove the following auxiliary results, which are used in the proof of Lemma D.4. Based on Lemmas D.3–D.5, we will show that all diagonal entries across all layers are identical, and likewise, all off-diagonal entries across layers are also equal.

Lemma D.3. Suppose we have a ground truth matrix $\mathbf{W}^* \in \mathbb{R}^{d \times d}$ whose diagonal entries are the same that we are observing, i.e., $w^* \triangleq w_{11}^* = w_{22}^* = \cdots = w_{dd}^*$ and $\Omega_{\mathrm{diag}}^{(d)} = \{(i,i)\}_{i=1}^d$. We factorize a solution matrix at time t as a product of L matrices,

$$W_{L:1}(t) = W_L(t)W_{L-1}(t)\cdots W_1(t), \quad W_l(t) \in \mathbb{R}^{d\times d} \quad \text{for all } l \in [L].$$

Suppose that for all $l \in [L]$ and $0 \le m \le k$, the following holds:

$$\boldsymbol{W}_{l}^{(m)}(t) = x^{(m)}\boldsymbol{I}_{d} + y^{(m)}\left(\boldsymbol{J}_{d} - \boldsymbol{I}_{d}\right),$$

for some scalars $x^{(m)}, y^{(m)} \in \mathbb{R}$ where we denote $\mathbf{A}^{(k)}(t)$ as k-th derivative with respect to t of a matrix $\mathbf{A}(t)$. Then, the k-th derivative of the product $\mathbf{W}_{L:1}(t)$ satisfies

$$w_{11}^{(k)}(t) = w_{22}^{(k)}(t) = \dots = w_{dd}^{(k)}(t).$$

Proof. Let us denote the m-th derivative of each layer matrix by

$$\boldsymbol{A}^{(m)} \triangleq \boldsymbol{W}_{l}^{(m)}(t).$$

Then, the k-th time derivative of the product $W_{L:1}(t)$ is given by the Leibniz rule:

$$\frac{d^k}{dt^k} \boldsymbol{W}_{L:1}(t) = \sum_{k_1 + \dots + k_L = k} {k \choose k_1, \dots, k_L} \boldsymbol{A}^{(k_L)} \boldsymbol{A}^{(k_{L-1})} \cdots \boldsymbol{A}^{(k_1)}.$$

By the assumption, each $A^{(m)}$ lies in the span of $\{I_d, J_d\}$, and since this span is closed under matrix multiplication and scalar multiplication (by Lemma D.2), each term in the sum lies in the same span. Hence, the entire sum $W^{(k)}(t)$ also lies in span $\{I_d, J_d\}$, which implies that all diagonal entries of $W^{(k)}(t)$ are equal.

Lemma D.4. Under the setting of Lemma D.3 where each factor matrix $W_l(0)$ is initialized according to (7), the following identities hold for all $k \in \mathbb{N} \cup \{0\}$ under the gradient flow dynamics defined in (3):

$$\left(\mathbf{W}_{l_1}^{(k)}(0) \right)_{ii} = \left(\mathbf{W}_{l_2}^{(k)}(0) \right)_{jj}, \quad i, j \in [d], \ l_1, l_2 \in [L],$$

$$\left(\mathbf{W}_{l_1}^{(k)}(0) \right)_{i_1 j_1} = \left(\mathbf{W}_{l_2}^{(k)}(0) \right)_{i_2 j_2}, \quad i_1 \neq j_1, i_2 \neq j_2 \in [d], l_1, l_2 \in [L].$$

Proof. For the base case, when k=0, these identities immediately follow from our initialization assumptions. Now, suppose the induction hypothesis holds for all orders m < k (with $k \ge 1$), which means we have:

$$\left(\boldsymbol{W}_{l_{1}}^{(m)}(0)\right)_{ii} = \left(\boldsymbol{W}_{l_{2}}^{(m)}(0)\right)_{jj}, \quad i, j \in [d], \ l_{1}, l_{2} \in [L],
\left(\boldsymbol{W}_{l_{1}}^{(m)}(0)\right)_{i_{1}j_{1}} = \left(\boldsymbol{W}_{l_{2}}^{(m)}(0)\right)_{i_{2}j_{2}}, \quad i_{1} \neq j_{1}, i_{2} \neq j_{2} \in [d], l_{1}, l_{2} \in [L].$$
(25)

By applying the Leibniz rule to (24), the k-th derivative of $W_l(t)$ is given by:

$$\boldsymbol{W}_{l}^{(k)}(t) = -\sum_{i_{1},\dots,i_{L}} {k-1 \choose i_{1},\dots,i_{L}} \prod_{r=l+1}^{L} \boldsymbol{W}_{r}^{(i_{r})}(t)^{\top} \cdot \nabla \ell(\boldsymbol{W}_{L:1}(t))^{(i_{l})} \cdot \prod_{r=1}^{l-1} \boldsymbol{W}_{r}^{(i_{r})}(t)^{\top}, \quad (26)$$

with $\sum_{l=1}^{L} i_l = k-1$ where each $i_l \geq 0$. Given our induction assumption in equation (25) for all m < k, let $x^{(m)}(0)$ denote the m-th derivative of the diagonal entries and $y^{(m)}(0)$ the m-th derivative of the off-diagonal entries at initialization. Note that at initialization, by Lemma D.3, under the assumption that $\boldsymbol{W}_l^{(m)}(0)$ lies in the span of $\{\boldsymbol{I}_d, \boldsymbol{J}_d\}$ leads to $w_{11}^{(m)}(0) = w_{22}^{(m)}(0) \cdots = w_{dd}^{(m)}(0)$. Therefore, we know $\nabla \ell(\boldsymbol{W}_{L:1}(0))^{(i_l)} = r^{(i_l)}(0)\boldsymbol{I}_d$ for all $i_l < k$, where $r^{(i_l)}(0) \triangleq r_{11}^{(i_l)}(0) = \cdots = r_{dd}^{(i_l)}(0)$. Thus, at initialization, since equation (26) consists of terms involving $x^{(m)}(0)$ and $y^{(m)}(0)$ for all m < k, we can rewrite the above expression at t = 0 in terms of these derivatives as follows:

$$W_{l}^{(k)}(0) = -\sum_{i_{1},...,i_{L}} {k-1 \choose i_{1},...,i_{L}} r^{(i_{l})}(0) \prod_{r \in [L] \setminus \{l\}} W_{r}^{(i_{r})}(0)$$

$$= -\sum_{i_{1},...,i_{L}} {k-1 \choose i_{1},...,i_{L}} r^{(i_{l})}(0) \prod_{r \in [L] \setminus \{l\}} (a_{r} \mathbf{I}_{d} + b_{r} \mathbf{J}_{d}),$$

where constants a_r and b_r are composed of $x^{(r)}(0)$ and $y^{(r)}(0)$. Then, by Lemma D.2, $W_l^{(k)}(0)$ can be expressed in terms of only two values—one for the diagonal entries and one for the off-diagonal entries:

$$W_l^{(k)}(0) = \alpha I_d + \beta J_d, \quad \alpha, \beta \in \mathbb{R},$$

thus concluding the proof.

 Lemma D.5. Under the setting of Lemma D.4, the symmetries are preserved for all time $t \ge 0$:

$$\begin{split} &(\boldsymbol{W}_{l_1}(t))_{ii} = (\boldsymbol{W}_{l_2}(t))_{jj} \quad \textit{for all } i, j \in [d], \ l_1, l_2 \in [L], \\ &(\boldsymbol{W}_{l_1}(t))_{i_1j_1} = (\boldsymbol{W}_{l_2}(t))_{i_2j_2} \quad \textit{for all } i_1 \neq j_1, i_2 \neq j_2 \in [d], \ l_1, l_2 \in [L]. \end{split}$$

Proof. By applying Lemma F.6 to the result of Lemma D.4, we can conclude that the symmetries are preserved for timesteps $t \ge 0$.

By the above lemmas, if the initialization follows the scheme in (7), then all diagonal entries of all layers are identical, and all off-diagonal entries are also identical. Under this condition, the gradient flow dynamics can be easily described by the following lemma.

Lemma D.6. Under the same conditions as in Lemma D.4, if the diagonal entries of each layer are identical at timestep t (denoted by x(t)), and if the off-diagonal entries of each layer are identical at timestep t (denoted by y(t)), then the time derivative of x(t) and y(t) are given as:

$$\dot{x}(t) = -\frac{(x(t) + (d-1)y(t))^{L-1} + (d-1)(x(t) - y(t))^{L-1}}{d}r(t),$$

$$\dot{y}(t) = -\frac{(x(t) + (d-1)y(t))^{L-1} - (x(t) - y(t))^{L-1}}{d}r(t).$$

Proof. For $l \in [L]$ the gradient flow dynamics of W_l are written as:

$$\dot{\mathbf{W}}_l(t) = -\prod_{i=l+1}^{L} \mathbf{W}_i(t)^{\top} \cdot \nabla \ell(\mathbf{W}_{L:1}(t)) \cdot \prod_{i=1}^{l-1} \mathbf{W}_i(t)^{\top},$$
(27)

where $\nabla \ell(\boldsymbol{W}_{L:1}(t)) = \operatorname{diag}(r(t), \dots, r(t))$. Since $\boldsymbol{W}_l(t)$ is comprised of x(t) in diagonal entries and y(t) in off-diagonal entries, the above dynamics can be rewritten as follows:

$$\dot{\mathbf{W}}_{l}(t) = -r(t) \left[\mathbf{W}_{l}(t) \right]^{L-l} \cdot \mathbf{I}_{d} \cdot \left[\mathbf{W}_{l}(t) \right]^{l-1}$$

$$= -r(t) \left[\mathbf{W}_{l}(t) \right]^{L-1}. \tag{28}$$

If we rewrite $W_l(t) = (x(t) - y(t))I_d + y(t)J_d$, its eigenvalues are derived as:

$$\lambda_1 = x(t) + (d-1)y(t)$$
 for the eigenvector $\mathbb{1}$, $\lambda_2 = x(t) - y(t)$ for any eigenvector orthogonal to $\mathbb{1}$ (multiplicity $d-1$).

 Here, we denote $\lambda_i \triangleq \lambda_i(\boldsymbol{W}_{L:1}(t))$, unless otherwise specified. Then, we can decompose $\boldsymbol{W}_l(t)$ with projection matrix $\boldsymbol{P}_{\parallel} = \frac{1}{d}\boldsymbol{J}_d$ and $\boldsymbol{P}_{\perp} = \boldsymbol{I}_d - \frac{1}{d}\boldsymbol{J}_d$ as follows:

$$\boldsymbol{W}_l(t) = \lambda_1 \boldsymbol{P}_{\parallel} + \lambda_2 \boldsymbol{P}_{\perp}.$$

Therefore, if we take (L-1)-th power of $W_l(t)$, we can derive:

$$\begin{aligned} [\boldsymbol{W}_{l}(t)]^{L-1} &= \lambda_{1}^{L-1} \boldsymbol{P}_{\parallel} + \lambda_{2}^{L-1} \boldsymbol{P}_{\perp} \\ &= \left(x(t) + (d-1)y(t) \right)^{L-1} \cdot \frac{1}{d} \boldsymbol{J}_{d} + \left(x(t) - y(t) \right)^{L-1} \left(\boldsymbol{I}_{d} - \frac{1}{d} \boldsymbol{J}_{d} \right) \\ &= \left(x(t) - y(t) \right)^{L-1} \boldsymbol{I}_{d} + \frac{\left(x(t) + (d-1)y(t) \right)^{L-1} - \left(x(t) - y(t) \right)^{L-1}}{d} \boldsymbol{J}_{d}. \end{aligned}$$

Recalling that I_d has 1 on the diagonal and 0 off-diagonal, and J_d has 1 in every entry, the entries of $[W_l(t)]^{L-1}$ are:

$$([\mathbf{W}_{l}(t)]^{L-1})_{ii} = (x(t) - y(t))^{L-1} + \frac{(x(t) + (d-1)y(t))^{L-1} - (x(t) - y(t))^{L-1}}{d}$$

$$= \frac{(x(t) + (d-1)y(t))^{L-1} + (d-1)(x(t) - y(t))^{L-1}}{d}, \quad \forall i \in [d],$$
 (29)

$$([\mathbf{W}_l(t)]^{L-1})_{ij} = \frac{(x(t) + (d-1)y(t))^{L-1} - (x(t) - y(t))^{L-1}}{d}, \quad \forall i \neq j \in [d].$$
(30)

This concludes the proof by substituting the above equations into equation (28). \Box

Under the gradient flow dynamics of the diagonal entry x(t) and y(t), we derive the dynamics of the singular value of $W_l(t)$.

Lemma D.7. Under the conditions of Lemma D.4, the singular values of $W_l(t)$, which is defined as $s_i(t)$ for $i \in [d]$, evolve according to:

$$\dot{s}_i(t) = -s_i^{L-1}(t)r(t), \quad i = 1, 2, \dots d.$$

Proof. By Lemma D.5, each factor matrix $W_l(t)$ is symmetric, having x(t) as its diagonal entries and y(t) as its off-diagonal entries. The distinct eigenvalues of $W_l(t)$ are $\lambda_1(t) = x(t) + (d-1)y(t)$ and $\lambda_2(t) = x(t) - y(t)$ (where $\lambda_2(t)$ has multiplicity d-1). Their time derivatives are calculated by:

$$\dot{\lambda_i}(t) = -\lambda_i^{L-1}(t)r(t),$$

Note that by setting m>1, we have $\lambda_1(0)\geq \lambda_2(0)>0$. If L=2, the solution of above equation is equal to $\lambda_i(t)=\lambda_i(0)\exp\left(-\int_0^t r(\tau)\mathrm{d}\tau\right)$, which means it maintains the positiveness of $\lambda_i(0)$ for all $t\geq 0$. For L>2, its general solution can be written as follows:

$$\lambda_i(t) = \left(\lambda_i(0)^{2-L} + (L-2) \int_0^t r(\tau) d\tau\right)^{\frac{1}{2-L}},$$

due to its positivity at initialization. Then, $\lambda_i(t)$ stays strictly positive, since it never reaches zero or changes sign. Therefore, due to the symmetry and positive definiteness of $W_l(t)$, we further conclude that $\lambda_i(t) \equiv s_i(t)$.

By the above lemma, we can solve the ODE and find $s_r(t)$ as follows:

$$s_r(t) = \begin{cases} s_r(0) \exp\left(-\int_0^t r(\tau) d\tau\right), & L = 2, \\ \left(s_r(0)^{2-L} + (L-2) \cdot \int_0^t r(\tau) d\tau\right)^{\frac{1}{2-L}}, & L > 2. \end{cases}$$

Since $s_1(0)=x(0)+(d-1)y(0)=\alpha\left(1+\frac{d-1}{m}\right)$ and $s_r(0)=x(0)-y(0)=\alpha(1-\frac{1}{m})$ for all $i\geq 2$, we can separate above equation as following:

$$s_{1}(t) = \begin{cases} \alpha \left(1 + \frac{d-1}{m}\right) \exp\left(-\int_{0}^{t} r(\tau) d\tau\right), & L = 2, \\ \left(\alpha^{2-L} \left(1 + \frac{d-1}{m}\right)^{2-L} + (L-2) \cdot \int_{0}^{t} r(\tau) d\tau\right)^{\frac{1}{2-L}}, & L > 2, \end{cases}$$

$$s_{r}(t) = \begin{cases} \alpha \left(1 - \frac{1}{m}\right) \exp\left(-\int_{0}^{t} r(\tau) d\tau\right), & L = 2, \\ \left(\alpha^{2-L} \left(1 - \frac{1}{m}\right)^{2-L} + (L-2) \cdot \int_{0}^{t} r(\tau) d\tau\right)^{\frac{1}{2-L}}, & L > 2. \end{cases}$$

$$r = 2, 3, \dots, d.$$

Then, we can establish a relationship between $s_1(t)$ and $s_r(t)$, thereby identifying an invariant property independent of time t:

• For
$$L=2$$
:
$$\frac{s_1(t)}{s_r(t)} = \frac{m+d-1}{m-1}, \tag{31}$$

• For L > 2:

$$s_1^{2-L}(t) - s_r^{2-L}(t) = \alpha^{2-L} \left(\left(1 + \frac{d-1}{m} \right)^{2-L} - \left(1 - \frac{1}{m} \right)^{2-L} \right). \tag{32}$$

Furthermore, we can derive a closed-form solution for the singular values by utilizing the convergence guarantee. From equation (29), the diagonal entries of the solution matrix can be expressed as:

$$w_{ii}(t) = ([\mathbf{W}_l(t)]^L)_{ii} = \frac{(x(t) + (d-1)y(t))^L + (d-1)(x(t) - y(t))^L}{d}, \quad \forall i \in [d].$$

Since $w_{ii}(t)$ converges to a fixed value w^* , and noting that s(t) = x(t) + (d-1)y(t) and $s_r(t) = x(t) - y(t)$, we obtain the following convergence equation:

$$w^* = \frac{s_1^L(\infty) + (d-1)s_r^L(\infty)}{d} = \frac{\sigma_1(\infty) + (d-1)\sigma_r(\infty)}{d},$$
 (33)

where we define $\sigma_i(t) \triangleq s_i^L(t)$ to denote the singular values of the product matrix, $W_{L:1}(t)$. Combining Equations (31) and (33), we derive a closed-form solution for the singular values of the depth-2 matrix as $t \to \infty$:

$$\sigma_1(\infty) = \left(\frac{w^*(m+d-1)^2}{m^2+d-1}\right)^{\frac{L}{2}},$$

$$\sigma_r(\infty) = \left(\frac{w^*(m-1)^2}{m^2+d-1}\right)^{\frac{L}{2}}, \quad r = 2, 3, \dots, d,$$

For the case when $L \geq 3$, we cannot obtain an exact analytical solution for $\sigma_r(\infty)$. Instead, we derive implicit equations for both $\sigma_1(\infty)$ and $\sigma_r(\infty)$ that cannot be easily solved without specifying numerical values:

$$\sigma_1^{\frac{2-L}{L}}(\infty) - \left(\frac{w^*d - \sigma_1(\infty)}{d-1}\right)^{\frac{2-L}{L}} = C_{\alpha,m,L,d},$$

$$(w^*d - (d-1)\sigma_r(\infty))^{\frac{2-L}{L}} - \sigma_r^{\frac{2-L}{L}}(\infty) = C_{\alpha,m,L,d}, \quad \text{for } r = 2, \dots, d.,$$

where $C_{\alpha,m,L,d} \triangleq \left(\frac{\alpha}{m}\right)^{2-L} \left(\left(m+d-1\right)^{2-L}-\left(m-1\right)^{2-L}\right)$. If we specify the values of $\alpha > 0, m > 1, d \geq 2, L \geq 3$ and $w^* > 0$ for ground-truth value, we can derive $\sigma_1(\infty)$ and $\sigma_r(\infty)$ of solution matrix of depth-L by substituting the values to above equations.

Remark. The $L\geq 3$ and $m=\infty$ case could arguably fall under the preceding analysis when other parameters are held fixed, as $m=\infty$ implies that all singular values are identical. However, a slight dependency on the specific value of α persists; for instance, tracking the overall result becomes challenging if α approaches zero while $m=\infty$. Therefore, we will restrict the scope of the aforementioned analysis to finite m. Consequently, the $L\geq 3$ and $m=\infty$ case will be analyzed separately in the following subsection.

D.3.2 Case for $L \geq 3$ and $m = \infty$

We now examine the $m=\infty$ case, which corresponds to an initialization scheme like $\mathbf{W}_l(0)=\alpha\mathbf{I}_d$. By Lemma D.1, the factor matrices $\mathbf{W}_l(t)$ remain diagonal for all $t\geq 0$, and thus the diagonal entries of the product matrix are $w_{ii}(t)=(\mathbf{W}_L(t))_{ii}(\mathbf{W}_{L-1}(t))_{ii}\cdots(\mathbf{W}_1(t))_{ii}$. Assuming zero-loss convergence is achieved for any initial choice of $\alpha>0$, it follows that $w_{ii}(\infty)=w^*$ for all i, and consequently, the overall matrix $\mathbf{W}_{L:1}(\infty)$ is diagonal with entries w^* .

Furthermore, let us consider the implications of Lemmas D.3–D.5. These lemmas hold under a condition y(t)=0, thereby belonging to span $\{I_d,J_d\}$, this leads to the result that each diagonal element of the factor matrices at convergence is $(W_l(\infty))_{ii}=(w^*)^{1/L}$ for all $i\in[d]$ and $l\in[L]$. This means each layer $W_l(\infty)$ becomes $(w^*)^{1/L}I_d$, and thus has identical singular values equal to $(w^*)^{1/L}$ (assuming $w^*\geq 0$). This, in turn, leads to the final claim that for the overall product matrix $W_{L:1}(\infty)$, its singular values $\sigma_i(\infty)$ satisfy $\sigma_i(\infty)=w^*$ for all $i\in[d]$.

D.3.3 Loss Convergence

We further establish loss convergence in the following proposition.

Proposition D.1. Let $W^* \in \mathbb{R}^{d \times d}$ be a ground-truth matrix with identical positive diagonal entries $w^* \triangleq w_{11}^* = \cdots = w_{dd}^* > 0$, and let $\Omega_{\mathrm{diag}}^{(d)} = \{(i,i)\}_{i=1}^d$. Consider gradient flow (3) on the product $W_{L:1}$, where each factor $W_l \in \mathbb{R}^{d \times d}$ is initialized as in (7). Define K from the initialization scale α by

$$K = \begin{cases} L\left(w_{ii}(0)\right)^{\frac{2L-2}{L}}, & 0 < w_{ii}(0) \le w^*, \\ L\left(w^*\right)^{\frac{2L-2}{L}}, & w_{ii}(0) \ge w^*, \end{cases}$$

where

$$w_{ii}(0) = \frac{\alpha^L \left((m+d-1)^L + (d-1)(m-1)^L \right)}{dm^L}.$$

Then, for all $t \geq 0$, the loss decays exponentially:

$$\ell(\mathbf{W}_{L:1}(t)) \le \ell(\mathbf{W}_{L:1}(0))e^{-2Kt}$$
.

Proof. Recall that the eigenvalues are given by $\lambda_1(t) = x(t) + (d-1)y(t)$ and $\lambda_2(t) = x(t) - y(t)$. From Lemma D.6, their time derivatives are

$$\dot{\lambda_1}(t) = -\lambda_1^{L-1}(t)r(t),$$

$$\dot{\lambda_2}(t) = -\lambda_2^{L-1}(t)r(t).$$

The diagonal entries $w_{ii}(t)$ of $W_{L:1}(t)$ can be written as

$$w_{ii}(t) = \frac{(x(t) + (d-1)y(t))^{L} + (d-1)(x(t) - y(t))^{L}}{d}$$
$$= \frac{\lambda_{1}^{L}(t) + (d-1)\lambda_{2}^{L}(t)}{d}.$$

Define the residual $r(t) = w_{ii}(t) - w^*$, where w^* is a constant. Differentiating r(t) and substituting the expressions for $\dot{\lambda}_1(t)$ and $\dot{\lambda}_2(t)$ yields

$$\dot{r}(t) = \frac{d}{dt} \left(w_{ii}(t) - w^* \right)
= \frac{L}{d} \lambda_1^{L-1}(t) \dot{\lambda}_1(t) + \frac{L(d-1)}{d} \lambda_2^{L-1}(t) \dot{\lambda}_2(t)
= \frac{L}{d} \lambda_1^{L-1}(t) \left(-\lambda_1^{L-1}(t)r(t) \right) + \frac{L(d-1)}{d} \lambda_2^{L-1}(t) \left(-\lambda_2^{L-1}(t)r(t) \right)
= -\left(\underbrace{\frac{L}{d} \lambda_1^{2L-2}(t) + \frac{L(d-1)}{d} \lambda_2^{2L-2}(t)}_{\triangleq K(t)} \right) r(t).$$
(34)

Thus $\dot{r}(t) = -K(t)r(t)$, whose solution is

$$r(t) = r(0) \exp\left(-\int_0^t K(\tau) d\tau\right). \tag{35}$$

Consequently, r(t) preserves the sign of r(0) for all $t \ge 0$. Also, by noting that the map $u \mapsto u^{\frac{2L-2}{L}}$ is convex on \mathbb{R}_+ , we can lower-bound K(t) using Jensen's inequality for any fixed t:

$$K(t) = L \left(\frac{\lambda_1^{2L-2}(t) + (d-1)\lambda_2^{2L-2}(t)}{d} \right)$$

$$= L \left(\frac{\left(\lambda_1^L(t)\right)^{\frac{2L-2}{L}} + (d-1)\left(\lambda_2^L(t)\right)^{\frac{2L-2}{L}}}{d} \right)$$

$$\geq L \left(\frac{\lambda_1^L(t) + (d-1)\lambda_2^L(t)}{d} \right)^{\frac{2L-2}{L}}$$

$$= L(w_{ii}(t))^{\frac{2L-2}{L}}. \tag{36}$$

Case 1 $(r(0) \le 0)$. Assume

$$0 < \alpha^{L} \le \frac{w^* dm^L}{(m+d-1)^L + (d-1)(m-1)^L},$$

which implies $r(0) \le 0$ and hence $r(t) \le 0$ by (35). For any $i \in \{1, 2\}$ with $\lambda_i(0) > 0$ we then have

$$\dot{\lambda}_i(t) = -\lambda_i^{L-1}(t)r(t) \ge 0,$$

so $\lambda_i(t) \ge \lambda_i(0) > 0$ for all $t \ge 0$, which in turn implies $w_{ii}(t) \ge w_{ii}(0)$. Therefore, we can lower bound (36) with $w_{ii}(0)$:

$$K(t) \ge L(w_{ii}(0))^{\frac{2L-2}{L}}.$$

Case 2 ($r(0) \ge 0$). If

$$\alpha^{L} \geq \frac{w^* dm^L}{(m+d-1)^L + (d-1)(m-1)^L},$$

then $r(0) \ge 0$ hence $r(t) \ge 0$ for all $t \ge 0$ by (35). Therefore, $w_{ii}(t) \ge w^*$, then we lower bound (36)

$$K(t) \ge L(w^*)^{\frac{2L-2}{L}}.$$

Moreover, since $\dot{\lambda}_i(t) = -\lambda_i^{L-1}(t)r(t) \le 0$, each $\lambda_i(t)$ is non-increasing. If it reaches 0 at some time, then $\dot{\lambda}_i(t) = 0$ there, so it cannot cross into the negative region; thus $\lambda_i(t) \ge 0$ for all $t \ge 0$. This justifies the use of (36).

By upper-bounding the absolute value of (35), we derive:

$$|r(t)| \le |r(0)| \exp(-Kt),$$

where $K=L(w_{ii}(0))^{\frac{2L-2}{L}}$ in Case 1 and $K=L(w^*)^{\frac{2L-2}{L}}$ in Case 2. Since $\ell(\boldsymbol{W}_{L:1}(t))=\frac{d}{2}r^2(t)$, we obtain the exponential decay of the loss:

$$\ell(\mathbf{W}_{L:1}(t)) \le \ell(\mathbf{W}_{L:1}(0)) \exp(-2Kt).$$

E Proof for Section 4

In this section, we provide the proofs for the propositions and theorems presented in Section 4. First, Subsection E.1 presents the general form of Proposition 4.1 along with its proof. Next, Subsection E.2 details the proof of Theorem 4.2, focusing on the 2×2 matrix case. Lastly, Subsection E.3 generalizes the core ideas of Theorem 4.2 to $d \times d$ matrices and provides the formal statement and the proof of Theorem 4.3.

E.1 GENERAL FORM AND PROOF OF PROPOSITION 4.1

We first present the general form of Proposition 4.1. This proposition applies to any "fully disconnected case", a scenario that involves the diagonal entries introduced within this same proposition.

For a $d \times d$ ground truth matrix W^* , the observed entries are given by $\Omega = \{(i_n, j_n)\}_{n=1}^d$. Since we consider the fully disconnected case, $i_n \neq i_m, j_n \neq j_m$ for all $n \neq m \in [d]$. We factorize the solution model at time t as $W_{A,B}(t) = A(t)B(t)$, where $W_{A,B}(t), A(t), B(t) \in \mathbb{R}^{d \times d}$. We consider the gradient flow dynamics with the loss function defined as in (2).

For a given row index k, since there exists a unique entry $(k,j) \in \Omega$, we denote this unique column index by $j^{(k)}$. Thus, $w_{k,j^{(k)}}^*$ and $w_{k,j^{(k)}}(t)$ refer to the ground truth weight $w_{k,j}^*$ and the time-varying weight $w_{k,j}(t)$ respectively, where $j=j^{(k)}$. Similarly, for a given column index l, since there exists a unique entry $(i,l) \in \Omega$, we denote this unique row index by $i^{(l)}$. Thus $w_{i^{(l)},l}^*$ and $w_{i^{(l)},l}$ refer to the ground truth weight $w_{i,l}^*$ and the time-varying weight $w_{i,l}(t)$ respectively, where $i=i^{(l)}$. Defining the residuals as $r_{ij}(t):=w_{ij}^*-w_{ij}(t)$, we adopt this compact notation for residuals as well. Then, we can derive a closed-form solution for arbitrary initialization with below proposition.

Proposition E.1. Consider a ground truth matrix $\mathbf{W}^* \in \mathbb{R}^{d \times d}$ and a set of d fully disconnected observations $\Omega = \{(i_n, j_n)\}_{n=1}^d$. The model is factorized as $\mathbf{W}_{A,\mathbf{B}}(t) = \mathbf{A}(t)\mathbf{B}(t)$, where the factors $\mathbf{A}(t), \mathbf{B}(t) \in \mathbb{R}^{d \times d}$. For each observed pair $(i_n, j_n) \in \Omega$, define the constants P_{i_n, j_n} and Q_{i_n, j_n} based on the initial values $\mathbf{A}(0)$ and $\mathbf{B}(0)$:

$$P_{i_n,j_n} \triangleq \sum_{k=1}^d a_{i_n,k}(0)b_{k,j_n}(0) \quad \text{and} \quad Q_{i_n,j_n} \triangleq \sum_{k=1}^d \left(a_{i_n,k}(0)^2 + b_{k,j_n}(0)^2\right).$$

Furthermore, for each such observed pair (i_n, j_n) , let the parameter \bar{r}_{i_n, j_n} be determined from the ground truth entry w_{i_n, j_n}^* and the constants defined above, as follows:

$$\bar{r}_{i_n,j_n} \triangleq \frac{1}{2} \log \left(\frac{P_{i_n,j_n} + \frac{Q_{i_n,j_n}}{2}}{w_{i_n,j_n}^* + \sqrt{w_{i_n,j_n}^*}^2 - P_{i_n,j_n}^2 + \left(\frac{Q_{i_n,j_n}}{2}\right)^2} \right).$$

Then, assuming convergence to a zero-loss solution (i.e., $w_{i_n,j_n}(\infty) = w_{i_n,j_n}^*$ for all $(i_n,j_n) \in \Omega$), any entry $a_{p,q}(\infty)$ of the converged matrix $A(\infty)$ and any entry $b_{p,q}(\infty)$ of the converged matrix $B(\infty)$ (for arbitrary indices $p,q \in [d]$) are explicitly given by:

$$\begin{split} a_{p,q}(\infty) &= a_{p,q}(0)\cosh\left(\bar{r}_{p,j^{(p)}}\right) - b_{q,j^{(p)}}(0)\sinh\left(\bar{r}_{p,j^{(p)}}\right), \\ b_{p,q}(\infty) &= b_{p,q}(0)\cosh\left(\bar{r}_{i^{(q)},q}\right) - a_{i^{(q)},p}(0)\sinh\left(\bar{r}_{i^{(q)},q}\right). \end{split}$$

Proof. We can express their evolution in the following vector form using the vectorized parameter $\boldsymbol{\theta}(t) \coloneqq \begin{bmatrix} \operatorname{vec}(\boldsymbol{A}(t)) \\ \operatorname{vec}(\boldsymbol{B}(t)) \end{bmatrix} \in \mathbb{R}^{2d^2}$:

$$\dot{\boldsymbol{\theta}}(t) = -\begin{bmatrix} \mathbf{0}_{d^2, d^2} & \boldsymbol{R}(t) \\ \boldsymbol{R}(t)^\top & \mathbf{0}_{d^2, d^2} \end{bmatrix} \boldsymbol{\theta}(t)$$
 (37)

where $\mathbf{R}(t) \in \mathbb{R}^{d^2 \times d^2}$ is defined as:

2162
2163
2164
2165
2166
2167
2168
2169
2170
2171 $R(t) = \begin{bmatrix} r_{1,j(1)}(t)e_{j(1)}^{\top} \\ r_{1,j(1)}(t)e_{j(1)+d}^{\top} \\ \vdots \\ r_{1,j(1)}(t)e_{j(1)+(d-1)d}^{\top} \\ r_{2,j(2)}(t)e_{j(2)}^{\top} \\ r_{2,j(2)}(t)e_{j(2)+d}^{\top} \\ \vdots \\ r_{2,j(2)}(t)e_{j(2)+d}^{\top} \end{bmatrix}$ (38)

for $e_i \in \mathbb{R}^{d^2}$ form the standard basis. Since $\begin{bmatrix} \mathbf{0}_{d^2,d^2} & \mathbf{R}(t) \\ \mathbf{R}(t)^\top & \mathbf{0}_{d^2,d^2} \end{bmatrix}$ commutes with any other t values, the solution is given as:

$$\boldsymbol{\theta}(t) = \exp\left(-\int_0^{\tau} \begin{bmatrix} \mathbf{0}_{d^2, d^2} & \boldsymbol{R}(t) \\ \boldsymbol{R}(t)^{\top} & \mathbf{0}_{d^2, d^2} \end{bmatrix} d\tau\right) \cdot \boldsymbol{\theta}(0)$$
 (39)

$$= \exp\left(-\begin{bmatrix} \mathbf{0}_{d^2,d^2} & \bar{\mathbf{R}}(t) \\ \bar{\mathbf{R}}(t)^\top & \mathbf{0}_{d^2 d^2} \end{bmatrix} d\tau\right) \cdot \boldsymbol{\theta}(0)$$
(40)

where

$$\bar{\boldsymbol{R}}(t) \coloneqq \int_{0}^{t} \boldsymbol{R}(\tau) \mathrm{d}\tau = \begin{bmatrix} \bar{r}_{1,j^{(1)}}(t)\boldsymbol{e}_{j^{(1)}}^{\top} \\ \bar{r}_{1,j^{(1)}}(t)\boldsymbol{e}_{j^{(1)}+d}^{\top} \\ \vdots \\ \bar{r}_{1,j^{(1)}}(t)\boldsymbol{e}_{j^{(1)}+(d-1)d}^{\top} \\ \bar{r}_{2,j^{(2)}}(t)\boldsymbol{e}_{j^{(2)}}^{\top} \\ \bar{r}_{2,j^{(2)}}(t)\boldsymbol{e}_{j^{(2)}+d}^{\top} \\ \vdots \\ \bar{r}_{d,j^{(d)}}(t)\boldsymbol{e}_{j^{(d)}+(d-1)d}^{\top} \end{bmatrix}$$

for $\bar{r}_{i,j}(t) = \int_0^t r_{i,j}(\tau) d\tau$. If we assume convergence, we get:

$$\boldsymbol{\theta}(\infty) = \exp\left(-\begin{bmatrix} \mathbf{0}_{d^2, d^2} & \bar{\boldsymbol{R}}(\infty) \\ \bar{\boldsymbol{R}}(\infty)^\top & \mathbf{0}_{d^2, d^2} \end{bmatrix} d\tau\right) \cdot \boldsymbol{\theta}(0)$$
(41)

$$= \left(\begin{bmatrix} \mathbf{I}_{d^2} & \mathbf{0}_{d^2,d^2} \\ \mathbf{0}_{d^2,d^2} & \mathbf{I}_{d^2} \end{bmatrix} - \begin{bmatrix} \mathbf{0}_{d^2,d^2} & \bar{\boldsymbol{R}}(t) \\ \bar{\boldsymbol{R}}(t)^{\top} & \mathbf{0}_{d^2,d^2} \end{bmatrix} + \frac{1}{2} \begin{bmatrix} \bar{\boldsymbol{R}}(t)\bar{\boldsymbol{R}}(t)^{\top} & \mathbf{0}_{d^2,d^2} \\ \mathbf{0}_{d^2,d^2} & \bar{\boldsymbol{R}}(t)^{\top}\bar{\boldsymbol{R}}(t) \end{bmatrix}$$
(42)

$$-\frac{1}{6} \begin{bmatrix} \mathbf{0}_{d^{2},d^{2}} & \bar{\mathbf{R}}(t)\bar{\mathbf{R}}(t)^{\top}\bar{\mathbf{R}}(t) \\ \bar{\mathbf{R}}(t)^{\top}\bar{\mathbf{R}}(t)\bar{\mathbf{R}}(t)^{\top} & \mathbf{0}_{d^{2},d^{2}} \end{bmatrix} + \frac{1}{24} \begin{bmatrix} (\bar{\mathbf{R}}(t)\bar{\mathbf{R}}(t)^{\top})^{2} & \mathbf{0}_{d^{2},d^{2}} \\ \mathbf{0}_{d^{2},d^{2}} & (\bar{\mathbf{R}}(t)^{\top}\bar{\mathbf{R}}(t))^{2} \end{bmatrix}$$
(43)

$$-\cdots$$
 $\cdot \boldsymbol{\theta}(0),$ (44)

which can be simplified as:

$$\boldsymbol{\theta}(\infty) = \begin{bmatrix} \boldsymbol{C} & \boldsymbol{D} \\ \boldsymbol{E} & \boldsymbol{F} \end{bmatrix} \boldsymbol{\theta}(0), \tag{45}$$

with C, D, E and F are defined as following:

$$\begin{split} & \boldsymbol{C} = \cosh \left(\mathrm{diag} \Big(\bar{r}_{1,j^{(1)}}, \ldots, \bar{r}_{1,j^{(1)}}, \bar{r}_{2,j^{(2)}}, \ldots, \bar{r}_{2,j^{(2)}}, \ldots, \bar{r}_{d,j^{(d)}}, \ldots, \bar{r}_{d,j^{(d)}} \Big) \right), \\ & \boldsymbol{F} = \cosh \left(\mathrm{diag} \Big(\bar{r}_{i^{(1)},1}, \bar{r}_{i^{(2)},2}, \ldots, \bar{r}_{i^{(d)},d}, \ldots, \bar{r}_{i^{(1)},1}, \bar{r}_{i^{(2)},2}, \ldots, \bar{r}_{i^{(d)},d} \Big) \right), \\ & \boldsymbol{D} = -\sinh \left(\Big[\bar{r}_{1,j^{(1)}} \boldsymbol{e}_{j^{(1)}}^{\top}, \ldots, \bar{r}_{1,j^{(1)}} \boldsymbol{e}_{j^{(1)}+(d-1)d}^{\top}, \ldots, \bar{r}_{d,j^{(d)}} \boldsymbol{e}_{j^{(d)}}^{\top}, \ldots, \bar{r}_{d,j^{(d)}} \boldsymbol{e}_{j^{(d)}+(d-1)d}^{\top} \Big]^{\top} \right), \\ & \boldsymbol{E} = -\sinh \left(\Big[\bar{r}_{1,j^{(1)}} \boldsymbol{e}_{j^{(1)}}, \ldots, \bar{r}_{1,j^{(1)}} \boldsymbol{e}_{j^{(1)}+(d-1)d}, \ldots, \bar{r}_{d,j^{(d)}} \boldsymbol{e}_{j^{(d)}}, \ldots, \bar{r}_{d,j^{(d)}} \boldsymbol{e}_{j^{(d)}+(d-1)d} \Big] \right). \end{split}$$

Here, for any matrix P, the operations $\cosh(P)$ and $\sinh(P)$ are performed elementwise. For a set of d observed indices Ω , there exists d corresponding unknown variables, \bar{r}_{i_k,j_k} . If convergence is guaranteed, the model yields d equations relating these variables to the d ground truth values. This implies that the variables \bar{r}_{i_k,j_k} can be characterized as a closed-form. To characterize more rigorously, we substitute C, D, E, and F into (45):

$$\boldsymbol{\theta}(\infty) = \begin{bmatrix} a_{1,1}(\infty) \\ a_{1,2}(\infty) \\ \vdots \\ a_{1,d}(\infty) \\ a_{2,1}(\infty) \\ a_{2,2}(\infty) \\ \vdots \\ a_{2,d}(\infty) \\ \vdots \\ a_{2,d}(\infty) \\ \vdots \\ a_{2,d}(\infty) \\ \vdots \\ a_{d,d}(\infty) \end{bmatrix} = \begin{bmatrix} a_{1,1}(0) \cosh(\bar{r}_{1,j^{(1)}}) - b_{1,j^{(1)}}(0) \sinh(\bar{r}_{1,j^{(1)}}) \\ a_{1,2}(0) \cosh(\bar{r}_{1,j^{(1)}}) - b_{2,j^{(1)}}(0) \sinh(\bar{r}_{1,j^{(1)}}) \\ \vdots \\ a_{1,d}(0) \cosh(\bar{r}_{1,j^{(1)}}) - b_{d,j^{(1)}}(0) \sinh(\bar{r}_{1,j^{(1)}}) \\ a_{2,1}(0) \cosh(\bar{r}_{2,j^{(2)}}) - b_{1,j^{(2)}}(0) \sinh(\bar{r}_{2,j^{(2)}}) \\ a_{2,2}(0) \cosh(\bar{r}_{2,j^{(2)}}) - b_{2,j^{(2)}}(0) \sinh(\bar{r}_{2,j^{(2)}}) \\ \vdots \\ a_{d,1}(\infty) \\ \vdots \\ a_{d,d}(\infty) \end{bmatrix} = \begin{bmatrix} a_{1,1}(0) \cosh(\bar{r}_{1,j^{(1)}}) - b_{1,j^{(1)}}(0) \sinh(\bar{r}_{1,j^{(1)}}) \\ a_{2,1}(0) \cosh(\bar{r}_{2,j^{(2)}}) - b_{2,j^{(2)}}(0) \sinh(\bar{r}_{2,j^{(2)}}) \\ \vdots \\ a_{2,d}(0) \cosh(\bar{r}_{2,j^{(2)}}) - b_{1,j^{(2)}}(0) \sinh(\bar{r}_{2,j^{(2)}}) \\ \vdots \\ a_{d,1}(0) \cosh(\bar{r}_{1,j^{(1)}}) - b_{1,j^{(d)}}(0) \sinh(\bar{r}_{2,j^{(2)}}) \\ \vdots \\ a_{d,1}(0) \cosh(\bar{r}_{1,j^{(1)}}) - b_{1,j^{(d)}}(0) \sinh(\bar{r}_{2,j^{(2)}}) \\ \vdots \\ a_{d,1}(0) \cosh(\bar{r}_{1,j^{(1)}}) - b_{1,j^{(2)}}(0) \sinh(\bar{r}_{2,j^{(2)}}) \\ \vdots \\ a_{d,1}(0) \cosh(\bar{r}_{1,j^{(2)}}) - b_{1,j^{(2)}}(0) \sinh(\bar{r}_{2,j^{(2)}}) \\ \vdots \\ a_{d,1}(0) \sinh(\bar{r}_{1,j^{(2)}}) - b_{1,j^{(2)}}(0) \sinh(\bar{r}_{2,j^{(2)}}) \\ \vdots \\ a_{d,1}(0) \cosh(\bar{r}_{1,j^{(2)}}) - b_{1,j^{(2)}}(0) \sinh(\bar{r}_{2,j^{(2)}}) \\ \vdots \\ a_{d,1}(0) \sinh(\bar{r}_{1,j^{(2)}}) - b_{1,j^{(2)}}(0)$$

Then, assuming convergence, for each observation $(i_n, j_n) \in \Omega$ (for n = 1, ..., d), we obtain the equation:

$$\begin{split} w_{i_n,j_n}^* &= w_{i_n,j_n}(\infty) = a_{i_n,1}(\infty)b_{1,j_n}(\infty) + \dots + a_{i_n,d}(\infty)b_{d,j_n}(\infty) \\ &= \sum_{k=1}^d \left[\left(a_{i_n,k}(0)\cosh(\bar{r}_{i_n,j_n}) - b_{k,j^{(i_n)}}(0)\sinh(\bar{r}_{i_n,j_n}) \right) \\ & \cdot \left(b_{k,j_n}(0)\cosh(\bar{r}_{i_n,j_n}) - a_{i_n,k}(0)\sinh(\bar{r}_{i_n,j_n}) \right) \right]. \end{split}$$

Let $C_n = \cosh(\bar{r}_{i_n,j_n})$ and $S_n = \sinh(\bar{r}_{i_n,j_n})$. Then we can rewrite the above equation as:

$$w_{i_{n},j_{n}}^{*} = \sum_{k=1}^{d} \left(a_{i_{n},k}(0) b_{k,j_{n}}(0) C_{n}^{2} - a_{i_{n},k}(0)^{2} C_{n} S_{n} - b_{k,j_{n}}(0)^{2} C_{n} S_{n} + a_{i_{n},k}(0) b_{k,j_{n}}(0) S_{n}^{2} \right)$$

$$= \left(\sum_{k=1}^{d} a_{i_{n},k}(0) b_{k,j_{n}}(0) \right) \left(C_{n}^{2} + S_{n}^{2} \right) - \left(\sum_{k=1}^{d} \left(a_{i_{n},k}(0)^{2} + b_{k,j_{n}}(0)^{2} \right) \right) C_{n} S_{n}$$

$$= P_{i_{n},j_{n}} \cosh(2\bar{r}_{i_{n},j_{n}}) - \frac{Q_{i_{n},j_{n}}}{2} \sinh(2\bar{r}_{i_{n},j_{n}}), \tag{47}$$

where $P_{i_n,j_n} = \sum_{k=1}^d a_{i_n,k}(0)b_{k,j_n}(0)$ and $Q_{i_n,j_n} = \sum_{k=1}^d \left(a_{i_n,k}(0)^2 + b_{k,j_n}(0)^2\right)$.

By solving (47) with respect to \bar{r}_{i_n,j_n} , we can get:

$$2w_{i_n,j_n}^* = P_{i_n,j_n} \left(e^{2\bar{r}_{i_n,j_n}} + e^{-2\bar{r}_{i_n,j_n}} \right) - \frac{Q_{i_n,j_n}}{2} \left(e^{2\bar{r}_{i_n,j_n}} - e^{-2\bar{r}_{i_n,j_n}} \right)$$
$$= e^{2\bar{r}_{i_n,j_n}} \left(P_{i_n,j_n} - \frac{Q_{i_n,j_n}}{2} \right) + e^{-2\bar{r}_{i_n,j_n}} \left(P_{i_n,j_n} + \frac{Q_{i_n,j_n}}{2} \right).$$

Multiply by $e^{2\bar{r}_{i_n,j_n}}$ leads to:

$$2w_{i_n,j_n}^* e^{2\bar{r}_{i_n,j_n}} = e^{4\bar{r}_{i_n,j_n}} \left(P_{i_n,j_n} - \frac{Q_{i_n,j_n}}{2} \right) + P_{i_n,j_n} + \frac{Q_{i_n,j_n}}{2}.$$

Rearrange into a quadratic equation by setting $u = e^{2\bar{r}_{i_n,j_n}}$:

$$\left(P_{i_n,j_n} - \frac{Q_{i_n,j_n}}{2}\right)u^2 - 2w_{i_n,j_n}^*u + P_{i_n,j_n} + \frac{Q_{i_n,j_n}}{2} = 0.$$

By solving the above equation while noting that $P_{i_n,j_n} - \frac{Q_{i_n,j_n}}{2} \le 0$ by the definition, we can get explicit solutions for \bar{r}_{i_n,j_n} :

$$\bar{r}_{i_n,j_n} = \frac{1}{2} \log \left(\frac{P_{i_n,j_n} + \frac{Q_{i_n,j_n}}{2}}{w_{i_n,j_n}^* + \sqrt{w_{i_n,j_n}^*}^2 - P_{i_n,j_n}^2 + \left(\frac{Q_{i_n,j_n}}{2}\right)^2} \right).$$

Note that each \bar{r}_{i_n,j_n} is solely determined by the initial points $\theta(0)$. With \bar{r}_{i_n,j_n} determined for each observed entry, we have closed-form expressions characterizing the model's learned relationship for these observations. Consequently, by (46), we have:

$$a_{p,q}(\infty) = a_{p,q}(0) \cosh\left(\bar{r}_{p,j^{(p)}}\right) - b_{q,j^{(p)}}(0) \sinh\left(\bar{r}_{p,j^{(p)}}\right),$$

$$b_{p,q}(\infty) = b_{p,q}(0) \cosh\left(\bar{r}_{i^{(q)},q}\right) - a_{i^{(q)},p}(0) \sinh\left(\bar{r}_{i^{(q)},q}\right).$$

E.2 PROOF OF THEOREM 4.2

In this section, we will provide the analysis of 2×2 matrix that starts from pre-trained weights with diagonal observations $w^* \triangleq w_{11}^* = w_{22}^*$, $W_{A,B}(t)$ cannot converge to a low-rank solution. Let $T_1 > t_1$ be the timestep that concludes the pre-train phase. For the sake of simplicity, we omit the ϵ term introduced in the pre-training phase. Then, we know from Proposition E.1, we have:

$$\mathbf{A}(T_1) = \mathbf{B}(T_1) = \begin{pmatrix} \sqrt{w^*} & 0\\ 0 & \sqrt{w^*} \end{pmatrix}. \tag{48}$$

In the post-train phase, we introduce an additional observation in the off-diagonal entries, specifically w_{12}^* or w_{21}^* . Without loss of generality, we assume $w_{12}^* > 0$ is revealed while other observations remain the same, i.e., $\Omega_{\text{post}} = \{(1,1),(1,2),(2,2)\}$. Note that the gradient of the post-train loss is:

$$\nabla \ell(\boldsymbol{W_{A,B}}) = \begin{pmatrix} w_{11} - w^* & w_{12} - w_{12}^* \\ 0 & w_{22} - w^* \end{pmatrix}$$
$$= \begin{pmatrix} a_{11}b_{11} + a_{12}b_{21} - w^* & a_{11}b_{12} + a_{12}b_{22} - w_{12}^* \\ 0 & a_{21}b_{12} + a_{22}b_{22} - w^* \end{pmatrix}.$$

For simplicity, we again omit the Ω term in the loss specification. We define the residuals for the relevant matrix elements as $r_{11} := w_{11} - w^*$, $r_{12} := w_{12} - w_{12}^*$, and $r_{22} := w_{22} - w^*$.

We begin by demonstrating a pairwise symmetry between the entries of $\boldsymbol{A}(t)$ and $\boldsymbol{B}(t)$, which simplifies subsequent analysis. To this end, we first provide the time derivatives for the elements of $\boldsymbol{A}(t)$ and $\boldsymbol{B}(t)$. Given the general gradient flow dynamics $\dot{\boldsymbol{A}}(t) = -\nabla \ell(\boldsymbol{W}_{\boldsymbol{A},\boldsymbol{B}}(t))\boldsymbol{B}^{\top}(t)$ and $\dot{\boldsymbol{B}}(t) = -\boldsymbol{A}^{\top}(t)\nabla \ell(\boldsymbol{W}_{\boldsymbol{A},\boldsymbol{B}}(t))$, the component-wise updates are as follows. For $\boldsymbol{A}(t)$:

$$\dot{a}_{11}(t) = b_{11}(t)(w^* - w_{11}(t)) + b_{12}(t)(w_{12}^* - w_{12}(t)),
\dot{a}_{12}(t) = b_{21}(t)(w^* - w_{11}(t)) + b_{22}(t)(w_{12}^* - w_{12}(t)),
\dot{a}_{21}(t) = b_{12}(t)(w^* - w_{22}(t)),
\dot{a}_{22}(t) = b_{22}(t)(w^* - w_{22}(t)),$$
(49)

and for $\boldsymbol{B}(t)$:

$$\dot{b}_{11}(t) = a_{11}(t)(w^* - w_{11}(t)),
\dot{b}_{12}(t) = a_{11}(t)(w_{12}^* - w_{12}(t)) + a_{21}(t)(w^* - w_{22}(t)),
\dot{b}_{21}(t) = a_{12}(t)(w^* - w_{11}(t)),
\dot{b}_{22}(t) = a_{12}(t)(w_{12}^* - w_{12}(t)) + a_{22}(t)(w^* - w_{22}(t)).$$
(50)

Using the equations above, we first present a result showing that the k-th derivative of each element in A(t) and B(t) at initialization exhibits a pairwise symmetry:

Lemma E.1. Let $W_{A,B}(T_1) = A(T_1)B(T_1) \in \mathbb{R}^{2 \times 2}$ be a product matrix, where $A(T_1)$ and $B(T_1)$ are matrices that are obtained at the end of the pre-training phase. Suppose the ground truth matrix satisfies $w_{11}^* = w_{22}^*$. Then for every $k \in \mathbb{N} \cup \{0\}$, the following identities hold:

$$a_{11}^{(k)}(T_1) = b_{22}^{(k)}(T_1), \quad a_{12}^{(k)}(T_1) = b_{12}^{(k)}(T_1), a_{21}^{(k)}(T_1) = b_{21}^{(k)}(T_1), \quad a_{22}^{(k)}(T_1) = b_{11}^{(k)}(T_1),$$

$$(51)$$

and consequently,

$$w_{11}^{(k)}(T_1) = w_{22}^{(k)}(T_1). (52)$$

Proof. We prove the statement by induction on k. When k = 0, by the initialization assumption, we have

$$a_{11}(T_1) = b_{22}(T_1), \quad a_{12}(T_1) = b_{12}(T_1), \quad a_{21}(T_1) = b_{21}(T_1), \quad a_{22}(T_1) = b_{11}(T_1),$$

2376 and therefore $w_{11}(T_1) = w_{22}(T_1)$.

Assume that for all orders m < k (with k > 1) the identities

$$a_{11}^{(m)}(T_1) = b_{22}^{(m)}(T_1), \quad a_{12}^{(m)}(T_1) = b_{12}^{(m)}(T_1), \quad a_{21}^{(m)}(T_1) = b_{21}^{(m)}(T_1), \quad a_{22}^{(m)}(T_1) = b_{11}^{(m)}(T_1),$$

hold, and hence also $w_{11}^{(m)}(T_1)=w_{22}^{(m)}(T_1)$. By the Leibniz rule, each element of the k-th derivative can be written as a finite sum involving derivatives of orders strictly less than k. For $\boldsymbol{A}(t)$:

$$\begin{split} a_{11}^{(k)}(t) &= -\sum_{j=0}^{k-1} \binom{k-1}{j} \left(b_{11}^{(k-1-j)}(t) r_{11}^{(j)}(t) + b_{12}^{(k-1-j)}(t) r_{12}^{(j)}(t) \right), \\ a_{12}^{(k)}(t) &= -\sum_{j=0}^{k-1} \binom{k-1}{j} \left(b_{21}^{(k-1-j)}(t) r_{11}^{(j)}(t) + b_{22}^{(k-1-j)}(t) r_{12}^{(j)}(t) \right), \\ a_{21}^{(k)}(t) &= -\sum_{j=0}^{k-1} \binom{k-1}{j} b_{12}^{(k-1-j)}(t) r_{22}^{(j)}(t), \\ a_{22}^{(k)}(t) &= -\sum_{j=0}^{k-1} \binom{k-1}{j} b_{22}^{(k-1-j)}(t) r_{22}^{(j)}(t), \end{split}$$

and for $\boldsymbol{B}(t)$:

$$\begin{split} b_{11}^{(k)}(t) &= -\sum_{j=0}^{k-1} \binom{k-1}{j} a_{11}^{(k-1-j)}(t) r_{11}^{(j)}(t), \\ b_{12}^{(k)}(t) &= -\sum_{j=0}^{k-1} \binom{k-1}{j} \left(a_{11}^{(k-1-j)}(t) r_{12}^{(j)}(t) + a_{21}^{(k-1-j)}(t) r_{22}^{(j)}(t) \right), \\ b_{21}^{(k)}(t) &= -\sum_{j=0}^{k-1} \binom{k-1}{j} a_{12}^{(k-1-j)}(t) r_{11}^{(j)}(t), \\ b_{22}^{(k)}(t) &= -\sum_{j=0}^{k-1} \binom{k-1}{j} \left(a_{12}^{(k-1-j)}(t) r_{12}^{(j)}(t) + a_{22}^{(k-1-j)}(t) r_{22}^{(j)}(t) \right). \end{split}$$

By the inductive hypothesis, all derivatives of order less than k satisfy the symmetric relations at $t=T_1$. Inserting these equalities into the expressions with $t=T_1$ above shows that the symmetry is maintained at the k-th order:

$$a_{11}^{(k)}(T_1) = b_{22}^{(k)}(T_1), \quad a_{12}^{(k)}(T_1) = b_{12}^{(k)}(T_1), \quad a_{21}^{(k)}(T_1) = b_{21}^{(k)}(T_1), \quad a_{22}^{(k)}(T_1) = b_{11}^{(k)}(T_1),$$
 proving equations (51) and (52).
$$\square$$

Lemma E.2. Under the setting of Lemma E.1, below relationships hold for all $t > T_1$:

$$a_{11}(t) = b_{22}(t), \quad a_{12}(t) = b_{12}(t),$$

 $a_{21}(t) = b_{21}(t), \quad a_{22}(t) = b_{11}(t),$

$$(53)$$

which further leads to $w_{11}(t) = w_{22}(t)$.

Proof. By Lemmas F.6 and E.1, we may conclude that for all $t \ge T_1$, equation (53) holds, and therefore $w_{11}(t) = w_{22}(t)$.

By Lemma E.2, all entries of B(t) can be expressed in terms of the entries of A(t) for all $t \ge T_1$. From this point onward, we will represent $W_{A,B}(t)$ solely using the elements of A(t). We begin by simplifying the time derivative of A(t) as follows:

2430
2431
$$\dot{a}_{11}(t) = a_{22}(t)(w^* - w_{11}(t)) + a_{12}(t)(w_{12}^* - w_{12}(t)),$$
2432
$$\dot{a}_{12}(t) = a_{21}(t)(w^* - w_{11}(t)) + a_{11}(t)(w_{12}^* - w_{12}(t)),$$
2433
$$\dot{a}_{21}(t) = a_{12}(t)(w^* - w_{22}(t)),$$
2434
$$\dot{a}_{22}(t) = a_{11}(t)(w^* - w_{22}(t)).$$
(54)

Rewriting $W_{A,B}(t)$ in terms of the elements of A(t) yields:

$$\mathbf{W}_{A,B}(t) = \mathbf{A}(t)\mathbf{B}(t)
= \begin{pmatrix} a_{11}(t) & a_{12}(t) \\ a_{21}(t) & a_{22}(t) \end{pmatrix} \begin{pmatrix} a_{22}(t) & a_{12}(t) \\ a_{21}(t) & a_{11}(t) \end{pmatrix}
= \begin{pmatrix} a_{11}(t)a_{22}(t) + a_{12}(t)a_{21}(t) & 2a_{11}(t)a_{12}(t) \\ 2a_{21}(t)a_{22}(t) & a_{11}(t)a_{22}(t) + a_{12}(t)a_{21}(t) \end{pmatrix}.$$
(55)

We can also simplify the time derivative of $W_{A,B}(t)$ as follows:

$$\dot{w}_{11}(t) = (w^* - w_{11}(t)) \left(a_{11}^2(t) + a_{12}^2(t) + a_{21}^2(t) + a_{22}^2(t) \right) + (w_{12}^* - w_{12}(t)) \left(a_{11}(t) a_{21}(t) + a_{12}(t) a_{22}(t) \right),$$

$$\dot{w}_{12}(t) = 2(w_{12}^* - w_{12}(t)) \left(a_{11}^2(t) + a_{12}^2(t) \right) + 2(w^* - w_{11}(t)) \left(a_{11}(t) a_{21}(t) + a_{12}(t) a_{22}(t) \right),$$

$$\dot{w}_{21}(t) = 2(w^* - w_{11}(t)) \left(a_{11}(t) a_{21}(t) + a_{12}(t) a_{22}(t) \right),$$

$$\dot{w}_{22}(t) = \dot{w}_{11}(t).$$
(56)

Using (55), we state the basic conservation law from Arora et al. (2018): if the matrices are initialized in a balanced manner, this balancedness is preserved throughout the training process. That is,

$$\boldsymbol{A}(T_1)^{\top}\boldsymbol{A}(T_1) = \boldsymbol{B}(T_1)\boldsymbol{B}(T_1)^{\top},$$

holds at initialization, this leads to

$$a_{11}^2(t) + a_{21}^2(t) = a_{12}^2(t) + a_{22}^2(t), \ \forall t \ge T_1.$$
 (57)

Now, we are going to examine the time derivative of the loss:

$$\frac{d}{dt}\ell(\boldsymbol{W}_{\boldsymbol{A},\boldsymbol{B}}(t)) = \left\langle \nabla \ell(\boldsymbol{W}_{\boldsymbol{A},\boldsymbol{B}}(t)), \dot{\boldsymbol{W}}(t) \right\rangle \\
= \left\langle \nabla \ell(\boldsymbol{W}_{\boldsymbol{A},\boldsymbol{B}}(t)), \dot{\boldsymbol{A}}(t)\boldsymbol{B}(t) + \boldsymbol{A}(t)\dot{\boldsymbol{B}}(t) \right\rangle \\
= \operatorname{Tr} \left(\nabla \ell^{\top}(\boldsymbol{W}_{\boldsymbol{A},\boldsymbol{B}}(t)) \left(\dot{\boldsymbol{A}}(t)\boldsymbol{B}(t) + \boldsymbol{A}(t)\dot{\boldsymbol{B}}(t) \right) \right) \\
= \operatorname{Tr} \left(\nabla \ell^{\top}(\boldsymbol{W}_{\boldsymbol{A},\boldsymbol{B}}(t)) \dot{\boldsymbol{A}}(t)\boldsymbol{B}(t) \right) + \operatorname{Tr} \left(\nabla \ell^{\top}(\boldsymbol{W}_{\boldsymbol{A},\boldsymbol{B}}(t))\boldsymbol{A}(t)\dot{\boldsymbol{B}}(t) \right) \\
= - \operatorname{Tr} \left(\nabla \ell^{\top}(\boldsymbol{W}_{\boldsymbol{A},\boldsymbol{B}}(t)) \nabla \ell(\boldsymbol{W}_{\boldsymbol{A},\boldsymbol{B}}(t)) \boldsymbol{B}^{\top}(t)\boldsymbol{B}(t) \right) \\
- \operatorname{Tr} \left(\nabla \ell^{\top}(\boldsymbol{W}_{\boldsymbol{A},\boldsymbol{B}}(t)) \boldsymbol{A}(t) \boldsymbol{A}^{\top}(t) \nabla \ell(\boldsymbol{W}_{\boldsymbol{A},\boldsymbol{B}}(t)) \right) \\
= - \operatorname{Tr} \left(\nabla \ell(\boldsymbol{W}_{\boldsymbol{A},\boldsymbol{B}}(t)) \boldsymbol{B}^{\top}(t) \boldsymbol{B}(t) \nabla \ell(\boldsymbol{W}_{\boldsymbol{A},\boldsymbol{B}}(t)) \right) \\
= - \operatorname{Tr} \left(\nabla \ell(\boldsymbol{W}_{\boldsymbol{A},\boldsymbol{B}}(t)) \boldsymbol{A}(t) \boldsymbol{A}^{\top}(t) \nabla \ell(\boldsymbol{W}_{\boldsymbol{A},\boldsymbol{B}}(t)) \right). \tag{58}$$

The third equality follows from the fact that for any two matrices A and B of the same size, $\langle A, B \rangle = \text{Tr}(A^{\top}B)$. The last equation holds due to the cyclic property of the trace. Combining (58) with Lemma F.7, we can ensure $L_1(t)$ and $L_2(t)$ are both positive semidefinite, which implies the loss is monotonically non-increasing for all $t \geq T_1$.

With Lemma E.2 and the monotonicity of the loss, we can guarantee positiveness of a_{11} , a_{22} , w_{11} , and w_{22} after the pre-train phase:

Lemma E.3. For a product matrix $W_{A,B}(t) = A(t)B(t) \in \mathbb{R}^{2\times 2}$, if $a_{11}(T_1), a_{22}(T_1), w_{11}(T_1),$ and $w_{22}(T_1)$ have all positive values, following inequalities hold for all $t \geq T_1$:

 $a_{11}(t), a_{22}(t) > 0, \quad a_{12}(t) \ge 0.$

Furthermore,

$$w_{11}(t), w_{22}(t) > 0$$

holds for all $t > T_1$.

Proof. We will prove the inequalities step by step.

Positiveness of a₁₁(t). For the sake of contradiction, assume that there exists a timestep $\tau_1 > T_1$ where $a_{11}(\tau_1) = 0$ holds. From (55) and Lemma F.3, we must have $\det(\mathbf{A}(\tau_1)) > 0$, which implies that $a_{12}(\tau_1)a_{21}(\tau_1) < 0$. Given the monotonicity of ℓ , $\mathbf{W}_{\mathbf{A},\mathbf{B}}(t)$ must satisfy:

$$\ell(\mathbf{W}_{\mathbf{A},\mathbf{B}}(t)) \le \ell(\mathbf{W}_{\mathbf{A},\mathbf{B}}(T_1)). \tag{59}$$

for all $t \geq T_1$. However, $W_{A,B}(\tau_1)$ cannot satisfy (59) because $w_{11}(\tau_1), w_{22}(\tau_1) < 0$ and $w_{12}(\tau_1) = 0$ for any $\tau_1 \geq 0$. This contradiction implies that such a τ_1 cannot exist.

Positiveness of a₂₂(t). Similarly, let's assume there exists a time $\tau_2 > T_1$ such that $a_{22}(\tau_2) = 0$ for the first time. We can express $W_{A,B}(\tau_2)$ as:

$$\boldsymbol{W_{A,B}}(\tau_2) = \begin{pmatrix} a_{12}(\tau_2)a_{21}(\tau_2) & 2a_{11}(\tau_2)a_{12}(\tau_2) \\ 0 & a_{12}(\tau_2)a_{21}(\tau_2) \end{pmatrix}.$$

where the diagonal entries are negative due to the condition $\det(\mathbf{A}(\tau_2)) > 0$. Therefore, the time derivative of a_{22} at timestep τ_2 is positive:

$$\dot{a}_{22}(\tau_2) = a_{11}(\tau_2)(w^* - w_{11}(\tau_2)) > 0.$$

Since $a_{22}(t)$ is increasing at point τ_2 , there exists time $t' < \tau_2$ such that $a_{22}(t') < 0$ (since $a_{22}(t)$ is continuous and differentiable), which is contradictory. Consequently, there cannot exist a τ_2 such that $a_{22}(\tau_2) = 0$.

Positiveness of a_{12}(t). Given that ℓ is non-decreasing, we can state:

$$\ell(\boldsymbol{W_{A,B}}(t)) = \frac{1}{2} \left[(w^* - w_{11}(t))^2 + (w_{12}^* - w_{12}(t))^2 + (w^* - w_{22}(t))^2 \right]$$

$$\leq \ell(\boldsymbol{W_{A,B}}(T_1)) = \frac{1}{2} w_{12}^{*2},$$

for all $t \ge T_1$. Since $(w^* - w_{11}(t))^2$ and $(w^* - w_{22}(t))^2$ are non-negative, $w_{12}(t)$ must be non-negative for all $t \ge T_1$. From (55), we know $w_{12}(t) = 2a_{11}(t)a_{12}(t)$, which implies $a_{12}(t) \ge 0$ for all $t \ge T_1$ with the above conclusion which states $a_{11}(t) > 0$.

Positiveness of w₁₁(t), w₂₂(t). Likewise, assume for the sake of contradiction that there exists a time $\tau_3 \ge T_1$ when $w_{11}(\tau_3) = 0$ is first satisfied. This directly implies that $a_{11}(\tau_3)a_{22}(\tau_3) = -a_{12}(\tau_3)a_{21}(\tau_3)$. Squaring both sides of the equation yields:

$$a_{11}^2(\tau_3)a_{22}^2(\tau_3) = a_{12}^2(\tau_3)a_{21}^2(\tau_3).$$

Subtracting $a_{12}^2(\tau_3)a_{22}^2(\tau_3)$ from both sides:

$$a_{11}^2(\tau_3)a_{22}^2(\tau_3) - a_{12}^2(\tau_3)a_{22}^2(\tau_3) = a_{12}^2(\tau_3)a_{21}^2(\tau_3) - a_{12}^2(\tau_3)a_{22}^2(\tau_3).$$

Factoring:

$$a_{22}^2(\tau_3)\left(a_{11}^2(\tau_3) - a_{12}^2(\tau_3)\right) = a_{12}^2(\tau_3)\left(a_{21}^2(\tau_3) - a_{22}^2(\tau_3)\right).$$

By the conservation law in (57), we have $a_{11}^2(\tau_3) + a_{21}^2(\tau_3) = a_{12}^2(\tau_3) + a_{22}^2(\tau_3)$, which leads to $a_{11}^2(\tau_3) - a_{12}^2(\tau_3) = a_{22}^2(\tau_3) - a_{21}^2(\tau_3)$. Replacing $a_{11}^2(\tau_3) - a_{12}^2(\tau_3)$ with $-(a_{21}^2(\tau_3) - a_{22}^2(\tau_3))$:

$$-a_{22}^2(\tau_3)\left(a_{21}^2(\tau_3)-a_{22}^2(\tau_3)\right)=a_{12}^2(\tau_3)\left(a_{21}^2(\tau_3)-a_{22}^2(\tau_3)\right).$$

This gives us:

$$\left(a_{12}^2(\tau_3) + a_{22}^2(\tau_3)\right) \left(a_{21}^2(\tau_3) - a_{22}^2(\tau_3)\right) = 0.$$

Since $a_{22}(\tau_3) > 0$ from the previous result, we can conclude that $a_{21}(\tau_3) = \pm a_{22}(\tau_3)$. To determine the sign of $a_{21}(\tau_3)$, recall that $W_{A,B}(\tau_3)$ is written as:

$$\mathbf{W}_{A,B}(au_3) = \begin{pmatrix} 0 & 2a_{11}(au_3)a_{12}(au_3) \\ 2a_{21}(au_3)a_{22}(au_3) & 0 \end{pmatrix}.$$

Since $a_{11}(\tau_3) > 0$, $a_{12}(\tau_3) \ge 0$ from the previous result, $2a_{11}(\tau_3)a_{12}(\tau_3) \ge 0$ holds. Also, given that $\det(\mathbf{W}_{\mathbf{A},\mathbf{B}}(\tau_3)) > 0$, we can determine that $a_{21}(\tau_3)$ is negative, which implies $a_{21}(\tau_3) = -a_{22}(\tau_3)$. Additionally, by the conservation law, we have $a_{11}^2(\tau_3) = a_{12}^2(\tau_3)$, which leads to $a_{11}(\tau_3) = a_{12}(\tau_3) > 0$.

Finally, consider the time derivative of w_{11} at timestep τ_3 , substituting $a_{11}(\tau_3)$ and $a_{21}(\tau_3)$ with $a_{12}(\tau_3)$ and $-a_{22}(\tau_3)$, respectively:

$$\dot{w}_{11}(\tau_3) = (w^* - w_{11}(\tau_3))(a_{11}^2(\tau_3) + a_{12}^2(\tau_3) + a_{21}^2(\tau_3) + a_{22}^2(\tau_3)) + (w_{12}^* - w_{12}(\tau_3))(a_{11}(\tau_3)a_{21}(\tau_3) + a_{12}(\tau_3)a_{22}(\tau_3)) = 2w^*(a_{12}^2(\tau_3) + a_{22}^2(\tau_3)) > 0,$$

which contradicts our initial assumption.

Given that the time derivative in the (56) includes the term $a_{11}(t)a_{21}(t) + a_{12}(t)a_{22}(t)$, we need to verify the sign of $a_{11}a_{21} + a_{12}a_{22}$ in order to proceed with the analysis. Below lemma shows that as long as $w_{12}(t) \le w_{12}^*$ holds, $a_{11}(t)a_{21}(t) + a_{12}(t)a_{22}(t)$ is always lower bounded by zero.

Lemma E.4. For a product matrix $W_{A,B}(t) = A(t)B(t) \in \mathbb{R}^{2\times 2}$, if at any point $t \in [T_1, T_2]$ we have $w_{12}(t) \leq w_{12}^*$, then the following inequality holds throughout the entire interval $[T_1, T_2]$:

$$a_{11}(t)a_{21}(t) + a_{12}(t)a_{22}(t) \ge 0.$$

Proof. We first define $g(t) \triangleq a_{11}(t)a_{21}(t) + a_{12}(t)a_{22}(t)$. Recall that at T_1 , we have $a_{12}(T_1) = a_{21}(T_1) = 0$, which implies $g(T_1) = 0$ as well. Note that by (54), at timestep T_1 , we have

$$\dot{a}_{12}(T_1) = a_{11}(T_1)(w_{12}^* - w_{12}(T_1)) + a_{21}(T_1)(w^* - w_{11}(T_1)) > 0,$$

while other elements remain unchanged. This indicates that g(t)>0 immediately after T_1 . We now show that if $g(\tau)>0$ for any $\tau\in (T_1,T_2]$, then there is no $\tau'\in [\tau,T_2]$ which satisfies both $g(\tau')=0$ and $\frac{d}{dt}g(t)\Big|_{t=\tau'}<0$. This implies that g(t) never becomes negative under the assumption of $w_{12}(t)\leq w_{12}^*$.

Suppose, for the sake of contradiction, that there exists a $\tau' \in [\tau, T_2]$ where $g(\tau') = 0$ and $\frac{d}{dt}g(t)\Big|_{t=\tau'} < 0$. Given $g(\tau') = 0$ and the conservation law in (57), and the inequalities from Lemma E.3, we can determine that there exist two combinations of the solution:

1.
$$a_{11}(\tau') = a_{22}(\tau'), \ a_{12}(\tau') = -a_{21}(\tau'), \ a_{11}(\tau') > a_{12}(\tau').$$

2.
$$a_{11}(\tau') = a_{22}(\tau'), \ a_{12}(\tau') = a_{21}(\tau') = 0.$$

We take the time derivative of g(t) at timestep τ' and substitute the values from (54) as follows:

$$\frac{d}{dt}g(t)\Big|_{t=\tau'} = \dot{a}_{11}(\tau')a_{21}(\tau') + a_{11}(\tau')\dot{a}_{21}(\tau') + \dot{a}_{12}(\tau')a_{22}(\tau') + a_{12}(\tau')\dot{a}_{22}(\tau')
= 2(w^* - w_{11}(\tau'))(a_{11}(\tau')a_{12}(\tau') + a_{21}(\tau')a_{22}(\tau'))
+ (w_{12}^* - w_{12}(\tau'))(a_{11}(\tau')a_{22}(\tau') + a_{12}(\tau')a_{21}(\tau')).$$
(60)

For the first case, substituting equations $a_{11}(\tau') = a_{22}(\tau')$ and $a_{12}(\tau') = -a_{21}(\tau')$ to (60) leads to:

$$\frac{d}{dt}g(t)\Big|_{t=\tau'} = (w_{12}^* - w_{12}(\tau'))w_{11}(\tau').$$

Since $w_{11}(t) > 0$ for all $t \ge T_1$, if $w_{12}(\tau') \le w_{12}^*$ holds, then g(t) cannot take negative values at time τ' .

For the second case, substituting equations $a_{11}(\tau') = a_{22}(\tau')$ and $a_{12}(\tau') = a_{21}(\tau') = 0$ to (60) leads to:

$$\frac{d}{dt}g(t)\Big|_{t=\tau'} = (w_{12}^* - w_{12}(\tau'))a_{11}^2(\tau'),$$

which is again a non-negative value if $w_{12}(\tau') \leq w_{12}^*$, leading to a contradiction.

Lemma E.5. For a product matrix $W_{A,B}(t) = A(t)B(t) \in \mathbb{R}^{2\times 2}$, the following inequalities holds for all timestep $t \geq T_1$:

$$w_{12}(t) \le w_{12}^*,$$

 $w_{11}(t), w_{22}(t) \ge w^*,$
 $w_{21}(t) \le 0.$

Proof. We will prove this lemma in several steps:

Step 1: $w_{12}(t) \leq w_{12}^*$ for all $t \geq T_1$.

We know $w_{12}(T_1) = 0 \le w_{12}^*$. Assume, for the sake of contradiction, that there exists a time $t' > T_1$ where t' is the first timestep such that $w_{12}(t') > w_{12}^*$. If this were true, there must exist a time s where $T_1 \le s < t'$ such that:

$$w_{12}(s) = w_{12}^*, \quad \dot{w}_{12}(s) > 0.$$

For these conditions to be met, $w_{12}(s)$ must satisfy:

$$\dot{w}_{12}(s) = 2(w^* - w_{11}(s))(a_{11}(s)a_{21}(s) + a_{12}(s)a_{22}(s)) > 0.$$
(61)

To satisfy (61), there are two possibilities:

$$(w^* - w_{11}(s)) > 0$$
 and $(a_{11}(s)a_{21}(s) + a_{12}(s)a_{22}(s)) > 0,$ (62)

or
$$(w^* - w_{11}(s)) < 0$$
 and $(a_{11}(s)a_{21}(s) + a_{12}(s)a_{22}(s)) < 0.$ (63)

However, neither of these can be true:

- 1. Equation (63) contradicts Lemma E.4, given that s < t'.
- 2. Equation (62) cannot be satisfied because there is no s where $w^* > w_{11}(s)$. If there were, there would be a time s' where $T_1 \le s' < s$ both satisfying $w_{11}(s') = w^*$, and $\dot{w}_{11}(s') < 0$. But we find:

$$\dot{w}_{11}(s') = (w_{12}^* - w_{12}(s'))(a_{11}(s')a_{21}(s') + a_{12}(s')a_{22}(s')) \ge 0.$$

This is because $w_{12}(s') < w_{12}^*$, and thus $a_{11}(s')a_{21}(s') + a_{12}(s')a_{22}(s') \ge 0$ by Lemma E.4. Therefore, our initial assumption must be false, implying that $w_{12}(t) \le w_{12}^*$ for all $t \ge T_1$.

Step 2: Prove $w_{11}(t) \ge w_{11}^*$ and $w_{22}(t) \ge w_{22}^*$ for all $t \ge T_1$.

Given $w_{12}(t) \le w_{12}^*$ for all $t \ge T_1$, Lemma E.4 implies $a_{11}(t)a_{21}(t) + a_{12}(t)a_{22}(t) \ge 0$ for all $t \ge T_1$. The evolution of w_{11} is given by:

$$\dot{w}_{11}(t) = (w^* - w_{11}(t))(a_{11}^2(t) + a_{12}^2(t) + a_{21}^2(t) + a_{22}^2(t)) + (w_{12}^* - w_{12}(t))(a_{11}(t)a_{21}(t) + a_{12}(t)a_{22}(t)).$$

By above equation, if there exists a time $t' \geq T_1$ where $w_{11}(t') = w^*$, we can conclude $\dot{w}_{11}(t') \geq 0$, and thus $w_{11}(t) \geq w^*$ for all $t \geq T_1$. By Lemma E.2, w_{22} has the same value as w_{11} , so $w_{22}(t) \geq w^*$ for all $t \geq T_1$.

Step 3: Prove $w_{21}(t) \leq 0$ for all $t \geq T_1$.

The evolution of w_{21} is given by:

$$\dot{w}_{21}(t) = 2(w^* - w_{11}(t))(a_{11}(t)a_{21}(t) + a_{12}(t)a_{22}(t)).$$

Since $w_{11}(t) \ge w^*$ and $a_{11}(t)a_{21}(t) + a_{12}(t)a_{22}(t) \ge 0$ for all $t \ge T_1$, we can conclude $w_{21}(t) \le 0$ for all $t \ge T_1$.

E.2.1 Proof of Loss Convergence

Recall that the time derivative of the loss function is written as:

$$\frac{d}{dt}\ell(\boldsymbol{W}_{\boldsymbol{A},\boldsymbol{B}}(t)) = -\operatorname{Tr}(\boldsymbol{L}_1(t)) - \operatorname{Tr}(\boldsymbol{L}_2(t)),$$

where $L_1(t)$ and $L_2(t)$ are defined in (58). To further our analysis, we can expand the time derivative of the loss by calculating the trace of $L_1(t)$ and $L_2(t)$. We omit the time index t when clear from context.

$$\begin{split} \boldsymbol{L}_{1} &= \begin{pmatrix} r_{11} & r_{12} \\ 0 & r_{22} \end{pmatrix} \begin{pmatrix} a_{21}^{2} + a_{22}^{2} & a_{11}a_{21} + a_{12}a_{22} \\ a_{11}a_{21} + a_{12}a_{22} & a_{11}^{2} + a_{12}^{2} \end{pmatrix} \begin{pmatrix} r_{11} & 0 \\ r_{12} & r_{22} \end{pmatrix} \\ &= \begin{pmatrix} r_{11}^{2}(a_{21}^{2} + a_{22}^{2}) + 2r_{11}r_{12}(a_{11}a_{21} + a_{12}a_{22}) + r_{12}^{2}(a_{11}^{2} + a_{12}^{2}) & C_{1} \\ C_{1} & r_{22}^{2}(a_{11}^{2} + a_{12}^{2}) \end{pmatrix}, \end{split}$$

for some time-dependent value C_1 . Following a similar process, we calculate L_2 :

$$L_{2} = \begin{pmatrix} r_{11} & 0 \\ r_{12} & r_{22} \end{pmatrix} \begin{pmatrix} a_{11}^{2} + a_{12}^{2} & a_{11}a_{21} + a_{12}a_{22} \\ a_{11}a_{21} + a_{12}a_{22} & a_{21}^{2} + a_{22}^{2} \end{pmatrix} \begin{pmatrix} r_{11} & r_{12} \\ 0 & r_{22} \end{pmatrix}$$

$$= \begin{pmatrix} r_{11}^{2}(a_{11}^{2} + a_{12}^{2}) & C_{2} \\ C_{2} & r_{12}^{2}(a_{11}^{2} + a_{12}^{2}) + 2r_{12}r_{22}(a_{11}a_{21} + a_{12}a_{22}) + r_{22}^{2}(a_{21}^{2} + a_{22}^{2}) \end{pmatrix},$$

again for the time-dependent value C_2 . With these expressions for L_1 and L_2 , we can now rewrite equation (58) in a more explicit form:

$$\frac{d}{dt}\ell(\boldsymbol{W}_{\boldsymbol{A},\boldsymbol{B}}(t)) = -\operatorname{Tr}(\boldsymbol{L}_{1}(t)) - \operatorname{Tr}(\boldsymbol{L}_{2}(t))$$

$$= -r_{11}^{2}(t) \left(a_{11}^{2}(t) + a_{12}^{2}(t) + a_{21}^{2}(t) + a_{22}^{2}(t)\right)$$

$$- 2r_{12}^{2}(t) \left(a_{11}^{2}(t) + a_{12}^{2}(t)\right)$$

$$- r_{22}^{2}(t) \left(a_{11}^{2}(t) + a_{12}^{2}(t) + a_{21}^{2}(t) + a_{22}^{2}(t)\right)$$

$$- 2r_{12}(t)r_{22}(t) \left(a_{11}(t)a_{21}(t) + a_{12}(t)a_{22}(t)\right)$$

$$- 2r_{11}(t)r_{12}(t) \left(a_{11}(t)a_{21}(t) + a_{12}(t)a_{22}(t)\right).$$
(64)

Note that the (64) is the non-positive term. Given that L_1 and L_2 are positive semi-definite, we can analyze each diagonal entry separately. This leads us to the following inequalities:

$$r_{11}^2(a_{21}^2 + a_{22}^2) + 2r_{11}r_{12}(a_{11}a_{21} + a_{12}a_{22}) + r_{12}^2(a_{11}^2 + b_{12}^2) \ge 0,$$

$$r_{12}^2(a_{11}^2 + a_{12}^2) + 2r_{12}r_{22}(a_{11}a_{21} + a_{12}a_{22}) + r_{22}^2(a_{21}^2 + a_{22}^2) \ge 0.$$

By rearranging the above inequalities, we obtain:

$$-2r_{11}r_{12}(a_{11}a_{21} + a_{12}a_{22}) \le r_{11}^2(a_{21}^2 + a_{22}^2) + r_{12}^2(a_{11}^2 + a_{12}^2),$$

$$-2r_{12}r_{22}(a_{11}a_{21} + a_{12}a_{22}) \le r_{12}^2(a_{11}^2 + a_{12}^2) + r_{22}^2(a_{21}^2 + a_{22}^2).$$

Substituting these inequalities into equation (64), we derive:

$$\frac{d}{dt}\ell(\mathbf{W}_{A,B}(t)) \le -r_{11}^2(t)\left(a_{11}^2(t) + a_{12}^2(t)\right) - r_{22}^2(t)\left(a_{11}^2(t) + a_{12}^2(t)\right). \tag{65}$$

This provides a tighter upper bound on the time derivative of the loss. However, it is still insufficient to guarantee convergence, as the bound does not depend on the term $r_{12}(t)$. As a result, even though the right-hand side converges to zero, this alone does not imply that the loss itself converges.

To further tighten the bound, we leverage the positive semidefiniteness of L_1 and L_2 . Specifically, note that for both QKQ^{\top} and $Q^{\top}KQ$ to be positive semi-definite, the only necessary condition is $K \succcurlyeq 0$. Therefore, we modify $L_1(t)$ to $\widetilde{L}_1(t) \triangleq \nabla \ell(W_{A,B}(t)) \left(B^{\top}(t)B(t) - \mu(t) \cdot e_2 e_2^{\top}\right) \nabla \ell^{\top}(W_{A,B}(t))$, where $\mu(t)$ is chosen to ensure that the matrix $B^{\top}(t)B(t) - \mu(t) \cdot e_2 e_2^{\top}$ remains positive semidefinite. This guarantees that $\widetilde{L}_1(t) \succcurlyeq 0$. To ensure this condition, $\mu(t)$ must satisfy:

$$\begin{aligned} \left| \boldsymbol{B}(t)^{\top} \boldsymbol{B}(t) - \mu(t) \cdot \boldsymbol{e}_{2} \boldsymbol{e}_{2}^{\top} \right) | &= \left| \begin{pmatrix} a_{21}^{2}(t) + a_{22}^{2}(t) & a_{11}(t) a_{21}(t) + a_{12}(t) a_{22}(t) \\ a_{11}(t) a_{21}(t) + a_{12}(t) a_{22}(t) & a_{11}^{2}(t) + a_{12}^{2}(t) - \mu(t) \end{pmatrix} \right| \\ &= - \left(a_{21}^{2}(t) + a_{22}^{2}(t) \right) \mu(t) + \left(a_{11}(t) a_{22}(t) - a_{12}(t) a_{21}(t) \right)^{2} \\ &> 0. \end{aligned}$$

Rearranging this inequality with respect to $\mu(t)$, we get:

$$\mu(t) \le \frac{(a_{11}(t)a_{22}(t) - a_{12}(t)a_{21}(t))^2}{a_{21}^2(t) + a_{22}^2(t)}$$

$$= \frac{\det(\boldsymbol{B}(t))^2}{a_{21}^2(t) + a_{22}^2(t)}.$$
(66)

Therefore, if we set $\mu(t)$ to satisfy the above inequality, we can guarantee \widetilde{L}_1 to be a positive semidefinite matrix. Now, $\widetilde{L}_1(t)$ can be calculated as:

$$\begin{split} \widetilde{\boldsymbol{L}_{1}} &= \begin{pmatrix} r_{11} & r_{12} \\ 0 & r_{22} \end{pmatrix} \begin{pmatrix} a_{21}^{2} + a_{22}^{2} & a_{11}a_{21} + a_{12}a_{22} \\ a_{11}a_{21} + a_{12}a_{22} & a_{11}^{2} + a_{12}^{2} - \mu \end{pmatrix} \begin{pmatrix} r_{11} & 0 \\ r_{12} & r_{22} \end{pmatrix} \\ &= \begin{pmatrix} r_{11}^{2}(a_{21}^{2} + a_{22}^{2}) + 2r_{11}r_{12}(a_{11}a_{21} + a_{12}a_{22}) + r_{12}^{2}(a_{11}^{2} + a_{12}^{2} - \mu) & \tilde{C} \\ \tilde{C} & r_{22}^{2}(a_{12}^{2} + a_{22}^{2} - \mu) \end{pmatrix}, \end{split}$$

for some \widetilde{C} . Since the matrix $\mathbf{B}^{\top}\mathbf{B} - \mu \cdot \mathbf{e}_2\mathbf{e}_2^{\top}$ is positive semi-definite, we can ensure $a_{12}^2 + a_{22}^2 - \mu \geq 0$. This leads to the following inequality from $\left(\widetilde{L_1}\right)_{11}$:

$$-2r_{11}r_{12}(a_{11}a_{21} + a_{12}a_{22}) \le r_{11}^2(a_{21}^2 + a_{22}^2) + r_{12}^2(a_{11}^2 + a_{12}^2 - \mu).$$

Finally, substituting this inequality into (64), we arrive at:

$$\frac{d}{dt}\ell(\mathbf{W}_{A,B}(t)) \le -\left(r_{11}^2(t) + r_{22}^2(t)\right)\left(a_{11}^2(t) + a_{12}^2(t)\right) - r_{12}^2(t)\mu(t). \tag{67}$$

To prove the convergence of the loss, our main remaining goal is to establish a time-invariant lower bound for

$$\min\left\{a_{11}^2(t)+a_{12}^2(t),\ \mu(t)\right\}$$

to apply Grönwall's inequality.

Lemma E.6. For a solution matrix $W_{A,B}(t)$ initialized as $W_{A,B}(T_1)$, which represents the state of the matrix after pre-training up to time T_1 , the inequality

$$\det\left(\boldsymbol{W_{A,B}}(t)\right) \ge w^{*2}$$

holds for all $t \geq T_1$.

Proof. Since $w_{12}(t)$ must satisfy $|w_{12}(t)-w_{12}^*| \leq \sqrt{2\ell(\pmb{W_{A,B}}(t))} \leq w_{12}^*$ by the monotonicity of the loss, we can ensure that $w_{12}(t) \geq 0$ for all $t \geq T_1$. Also, by Lemma E.5, we have $w_{11}(t), w_{22}(t) \geq w^*$, and $w_{21}(t) \leq 0$ for all $t \geq T_1$. Under these conditions, $\det(\pmb{W_{A,B}}(t))$ can be lower bounded as:

$$\det(\mathbf{W}_{A,B}(t)) = w_{11}(t)w_{22}(t) - w_{12}(t)w_{21}(t) \ge w^{*2},$$

for all timesteps $t \geq T_1$.

Lemma E.7. For $\mu(t)$ defined to satisfy (66) and the entries in A(t), the following inequality holds for all timesteps $t \ge T_1$:

$$\min\left\{a_{11}^2(t) + a_{12}^2(t), \ \mu(t)\right\} \ge w^*.$$

Proof. To prove the lower bound of $a_{11}^2(t) + a_{12}^2(t)$, Our goal is to demonstrate that $a_{11}^2(t) + a_{12}^2(t) \ge w^*$ for all timesteps t after T_1 . By Lemma E.7, we have $\|\mathbf{W}_{\mathbf{A},\mathbf{B}}(t)\|_F \ge \sqrt{2}w^*$, which leads to:

$$\begin{split} \sqrt{2}w^* &\leq \|\boldsymbol{W}_{\boldsymbol{A},\boldsymbol{B}}(t)\|_F \\ &= \sqrt{\sigma_1^2\left(\boldsymbol{W}_{\boldsymbol{A},\boldsymbol{B}}(t)\right) + \sigma_2^2\left(\boldsymbol{W}_{\boldsymbol{A},\boldsymbol{B}}(t)\right)}. \end{split}$$

By applying Lemma F.4, we have:

$$\sqrt{\sigma_1^2(\boldsymbol{W}_{\boldsymbol{A},\boldsymbol{B}}(t)) + \sigma_2^2(\boldsymbol{W}_{\boldsymbol{A},\boldsymbol{B}}(t))} = \sqrt{\sigma_1^4(\boldsymbol{A}(t)) + \sigma_2^4(\boldsymbol{A}(t))}$$

$$= \sqrt{(\sigma_1^2(\boldsymbol{A}(t)) + \sigma_2^2(\boldsymbol{A}(t)))^2 - 2\sigma_1^2(\boldsymbol{A}(t))\sigma_2^2(\boldsymbol{A}(t))}$$

$$= \sqrt{\|\boldsymbol{A}(t)\|_F^4 - 2\det(\boldsymbol{A}(t))^2}.$$
(68)

Rewriting (68) while applying Lemmas F.4 and E.6 leads to:

$$\|\mathbf{A}(t)\|_F^4 \ge 2w^{*2} + 2\det(\mathbf{A}(t))^2$$

= $2w^{*2} + 2\det(\mathbf{W}_{\mathbf{A},\mathbf{B}}(t))$
> $4w^{*2}$.

Thus, $\boldsymbol{A}(t)$ have to satisfy $\|\boldsymbol{A}(t)\|_F^2 \geq 2w^*$ for all timesteps $t \geq T_1$. Now, assume that there exists a time $t' > T_1$ such that $a_{11}^2(t') + a_{12}^2(t') < w^*$. To satisfy inequality $\|\boldsymbol{A}(t')\|_F^2 \geq 2w^*$, we would need at least $a_{21}^2(t') + a_{22}^2(t') > w^*$ to hold. To verify the value of $a_{21}^2(t') + a_{22}^2(t')$, we take its time derivative using (54):

$$\begin{split} \frac{d}{dt}(a_{21}^2(t) + a_{22}^2(t)) &= 2a_{21}(t)a_{21}(t) + 2a_{22}(t)a_{22}(t) \\ &= -2a_{12}(t)a_{21}(t)r_{22}(t) - 2a_{11}(t)a_{22}(t)r_{22}(t) \\ &= -2r_{22}(t)(a_{11}(t)a_{22}(t) + a_{12}(t)a_{21}(t)) \\ &= 2w_{11}(t)(w^* - w_{11}(t)). \end{split}$$

Since $w_{11}(t) \geq w^*$ holds by Lemma E.5 for all $t \geq T_1$, we conclude $a_{21}^2(t) + a_{22}^2(t)$ is monotonically non-increasing from time $t \geq T_1$. Since $a_{12}^2(T_1) + a_{22}^2(T_1)$ is initialized as w^* , this implies that $a_{21}^2(t') + a_{22}^2(t') \leq w^*$. Consequently, there cannot exist a $t' > T_1$ such that $a_{11}^2(t') + a_{12}^2(t') < w^*$ holds, which leads to contradiction.

Next, we are now showing that the term $\frac{\det(B(t))^2}{a_{21}^2(t)+a_{22}^2(t)}$ is lower bounded by w^* . Therefore, if we set $\mu(t)$ as w^* , we can guarantee the positive semidefiniteness of $\widetilde{L}_1(t)$.

By applying Lemma F.4 and the lower bound of $\det(W_{A,B}(t))$ by Lemma E.6, we have

$$\frac{\det (\boldsymbol{B}(t))^2}{a_{21}^2(t) + a_{22}^2(t)} = \frac{\det (\boldsymbol{W}_{\boldsymbol{A},\boldsymbol{B}}(t))}{a_{21}^2(t) + a_{22}^2(t)} \ge \frac{w^{*2}}{a_{21}^2(t) + a_{22}^2(t)}$$

Also, from the previous result, we have an upper bound on $a_{21}^2(t) + a_{22}^2(t)$, which is $a_{21}^2(t) + a_{22}^2(t) \le w^*$. Combining these results, the following inequality holds:

$$\frac{\det{(\pmb{W_{A,B}}(t))}}{a_{21}^2(t) + a_{22}^2(t)} \ge w^*.$$

Therefore, if we set $\mu(t)$ to be w^* , $\mu(t)$ can satisfy the positive semidefiniteness condition. By combining the results, we can finally guarantee:

$$\min\left\{a_{11}^2(t)+a_{12}^2(t),\;\mu(t)\right\}\geq w^*.$$

Using the results of Lemma E.7, we can rewrite (67) as follows:

$$\frac{d}{dt}\ell(\mathbf{W}_{A,B}(t)) \le -\left(r_{11}^2(t) + r_{22}^2(t)\right)\left(a_{11}^2(t) + a_{12}^2(t)\right) - r_{12}^2(t)\mu(t)
\le -\left(r_{11}^2(t) + r_{12}^2(t) + r_{22}^2(t)\right)w^*
\le -2w^*\ell(\mathbf{W}_{A,B}(t)).$$

Applying Grönwall's inequality to our previous result, we can now demonstrate loss convergence where $t \ge T_1$:

$$\ell(\mathbf{W}_{A,B}(t)) \le \ell(\mathbf{W}_{A,B}(T_1))e^{-2w^*(t-T_1)}$$

$$= \frac{1}{2}w_{12}^{*2}e^{-2w^*(t-T_1)}.$$
(69)

This inequality allows us to conclude that $\ell(W_{A,B}(t))$ converges to zero exponentially.

E.2.2 PROOF OF STABLE RANK BOUND

From (69), we know that at convergence, $w_{11}(\infty) = w_{22}(\infty) = w^*$ and $w_{12}(\infty) = w^*_{12}$. Although a closed-form expression for $w_{21}(\infty)$ is unavailable, Lemma E.5 shows that $w_{21}(t) \le 0$ for $t \ge T_1$, which implies $w_{21}(\infty) \le 0$. This indicates that the test loss remains strictly positive, as the ground-truth value $w^*_{21} = \frac{w^{*2}}{w^*_{12}}$ is assumed to be strictly positive.

In this section, we leverage the fast convergence rate detailed in (69) to establish bounds on the singular values of the converged matrix $W_{A,B}(\infty)$. Subsequently, these singular value bounds are used to further bound the stable rank of $W_{A,B}(\infty)$.

Lemma E.8. The singular values of $W_{A,B}(\infty)$ fulfill:

$$\sigma_1(\boldsymbol{W}_{\boldsymbol{A},\boldsymbol{B}}(\infty)) \le w^* \cdot \exp\left(2\frac{w_{12}^*}{w^*}\right),$$

$$\sigma_2(\boldsymbol{W}_{\boldsymbol{A},\boldsymbol{B}}(\infty)) \ge w^* \cdot \exp\left(-2\frac{w_{12}^*}{w^*}\right).$$

Proof. We denote the singular values of $W_{A,B}(t)$ as $\sigma_r(t)$ for simplicity. By Lemma F.1, we can get general solution of each singular value $\sigma_r(t)$ by solving linear differential equation:

$$\sigma_r(t) = \sigma_r(s) \cdot \exp\left(-2\int_{t'=s}^t \langle \nabla \ell(\boldsymbol{W}_{\boldsymbol{A},\boldsymbol{B}}(t')), \boldsymbol{u}_r(t')\boldsymbol{v}_r^{\top}(t')\rangle dt'\right), \quad r = 1, 2, \quad (70)$$

where $u_r(t)$ and $v_r(t)$ denotes left and right singular vector of corresponding r-th singular value, respectively. Since $u_r(t)$ and $v_r(t)$ are both unit vectors, applying Cauchy-Schwartz inequality, we can bound $\langle \nabla \ell(W_{A,B}(t)), u_r(t)v_r^{\top}(t) \rangle$ by:

$$\begin{split} \left| \left\langle \nabla \ell(\boldsymbol{W}_{\boldsymbol{A},\boldsymbol{B}}(t)), \boldsymbol{u}_r(t) \boldsymbol{v}_r^\top(t) \right\rangle \right| &\leq \left\| \nabla \ell(\boldsymbol{W}_{\boldsymbol{A},\boldsymbol{B}}(t)) \right\|_F \cdot \left\| \boldsymbol{u}_r(t) \boldsymbol{v}_r^\top(t) \right\|_F \\ &= \left\| \nabla \ell(\boldsymbol{W}_{\boldsymbol{A},\boldsymbol{B}}(t)) \right\|_F \\ &= \sqrt{2\ell(\boldsymbol{W}_{\boldsymbol{A},\boldsymbol{B}}(t))}. \end{split}$$

we can get bound $\sigma_r(t)$ as following:

$$\sigma_r(s) \cdot \exp\left(-2\sqrt{2} \int_{t'=s}^t \sqrt{\ell(\boldsymbol{W}_{\boldsymbol{A},\boldsymbol{B}}(t'))} dt'\right) \le \sigma_r(t) \le \sigma_r(s) \cdot \exp\left(2\sqrt{2} \int_{t'=s}^t \sqrt{\ell(\boldsymbol{W}_{\boldsymbol{A},\boldsymbol{B}}(t'))} dt'\right)$$
(71)

With the setting above, in the pre-train section, after T_1 timesteps, we prove that $\sigma_1(T_1) = \sigma_2(T_1) = w^*$. Starting from T_1 with pre-trained weights, we can lower bound $\sigma_2(\mathbf{W}_{A,B}(t))$ with equations (69)

and (71) when $t \ge T_1$ as follows:

$$\sigma_{2}(t) \geq \sigma_{2}(T_{1}) \cdot \exp\left(-2\sqrt{2} \int_{t'=T_{1}}^{t} \sqrt{\ell(\mathbf{W}_{\mathbf{A},\mathbf{B}}(t'))} dt'\right)$$

$$\geq w^{*} \cdot \exp\left(-2w_{12}^{*} \int_{t'=T_{1}}^{t} e^{-w^{*}(t'-T_{1})} dt'\right)$$

$$= w^{*} \cdot \exp\left(-\frac{2w_{12}^{*}}{w^{*}} \left(1 - e^{-w^{*}(t-T_{1})}\right)\right).$$

and when $t \to \infty$, $\sigma_2(\infty)$ can be lower bounded by:

$$\sigma_2(\infty) \ge w^* \cdot e^{-2 \cdot \frac{w_{12}^*}{w^*}}.$$

In the same way, we can upper bound $\sigma_1(\infty)$ by:

$$\sigma_1(\infty) \le w^* \cdot e^{2 \cdot \frac{w_{12}^*}{w^*}}.$$

By Lemma E.8, we can now lower bound the stable rank of a matrix $W_{A,B}(\infty)$:

$$\begin{split} \frac{\|\boldsymbol{W}_{\boldsymbol{A},\boldsymbol{B}}(\infty)\|_F^2}{\|\boldsymbol{W}_{\boldsymbol{A},\boldsymbol{B}}(\infty)\|_2^2} &= \frac{\sigma_1^2(\boldsymbol{W}_{\boldsymbol{A},\boldsymbol{B}}(\infty)) + \sigma_2^2(\boldsymbol{W}_{\boldsymbol{A},\boldsymbol{B}}(\infty))}{\sigma_1^2(\boldsymbol{W}_{\boldsymbol{A},\boldsymbol{B}}(\infty))} \\ &= 1 + \frac{\sigma_2^2(\boldsymbol{W}_{\boldsymbol{A},\boldsymbol{B}}(\infty))}{\sigma_1^2(\boldsymbol{W}_{\boldsymbol{A},\boldsymbol{B}}(\infty))} \\ &\geq 1 + \exp\left(-8\frac{w_{12}^*}{w^*}\right), \end{split}$$

which concludes the proof of Theorem 4.2.

E.3 FORMAL STATEMENT AND PROOF OF THEOREM 4.3

We now extend the preceding analysis to the general case involving a ground truth matrix $W^* \in \mathbb{R}^{d \times d}$. The solution matrix $W_{A,B} \in \mathbb{R}^{d \times d}$ is again factorized as $W_{A,B} = AB$, where both $A, B \in \mathbb{R}^{d \times d}$. In this section, our detailed presentation and proof of Theorem 4.3 (from the main text) are structured as follows: we first introduce and prove Theorem E.2, which is then followed by its direct consequence, Corollary E.3.

We use the slightly modified loss function:

$$\mathcal{L}(\boldsymbol{A}, \boldsymbol{B}) = \frac{1}{2} \sum_{n=1}^{N} \left(\langle \boldsymbol{A} \boldsymbol{B}, \boldsymbol{X}_{n} \rangle - y_{n} \right)^{2}, \tag{72}$$

where the measurement matrix $\boldsymbol{X}_n = \boldsymbol{e}_{i_n} \boldsymbol{e}_{j_n}^{\top}$ represents a masking matrix, with the n-th observed entry set to one and all other entries set to zero, and $y_n \in \mathbb{R}$ denotes the ground truth value of the n-th observation. Then, by defining $\boldsymbol{\Theta} = \begin{bmatrix} \boldsymbol{A} \\ \boldsymbol{B}^{\top} \end{bmatrix} \in \mathbb{R}^{2d \times d}$ and $\bar{\boldsymbol{X}}_n = \frac{1}{2} \begin{bmatrix} \boldsymbol{0} & \boldsymbol{X}_n \\ \boldsymbol{X}_n^{\top} & \boldsymbol{0} \end{bmatrix} \in \mathbb{R}^{2d \times 2d}$, we can rewrite the (72) as:

$$\mathcal{L}(\boldsymbol{A}, \boldsymbol{B}) = \tilde{\mathcal{L}}(\boldsymbol{\Theta}) = \frac{1}{2} \sum_{n=1}^{N} \left(\langle \boldsymbol{\Theta} \boldsymbol{\Theta}^{\top}, \bar{\boldsymbol{X}}_{n} \rangle - y_{n} \right)^{2}$$
$$= \frac{1}{2} \| F(\boldsymbol{\Theta}) - \boldsymbol{y} \|_{2}^{2}. \tag{73}$$

Here, $F(\Theta)$ and y represent vectors defined as:

$$F(\mathbf{\Theta}) \triangleq \begin{bmatrix} \langle \mathbf{\Theta} \mathbf{\Theta}^{\top}, \bar{\mathbf{X}}_{1} \rangle \\ \langle \mathbf{\Theta} \mathbf{\Theta}^{\top}, \bar{\mathbf{X}}_{2} \rangle \\ \vdots \\ \langle \mathbf{\Theta} \mathbf{\Theta}^{\top}, \bar{\mathbf{X}}_{N} \rangle \end{bmatrix} \in \mathbb{R}^{N}, \quad \mathbf{y} \triangleq \begin{bmatrix} y_{1} \\ y_{2} \\ \vdots \\ y_{N} \end{bmatrix} \in \mathbb{R}^{N}.$$
(74)

By reparameterizing A, B to Θ , and X_n to \bar{X}_n , we can reduce the parameter matrices into a single matrix Θ while ensuring the symmetry of $\Theta\Theta^{\top}$. We train the model Θ via gradient flow, where the loss evolution is given by:

$$\dot{\bar{\mathcal{L}}}(\boldsymbol{\Theta}(t)) = (F(\boldsymbol{\Theta}(t)) - \boldsymbol{y})^{\top} \dot{F}(\boldsymbol{\Theta}(t))$$

$$= (F(\boldsymbol{\Theta}(t)) - \boldsymbol{y})^{\top} \begin{bmatrix}
\frac{d}{dt} \langle \boldsymbol{\Theta}(t) \boldsymbol{\Theta}(t)^{\top}, \bar{\boldsymbol{X}}_{1} \rangle \\
\frac{d}{dt} \langle \boldsymbol{\Theta}(t) \boldsymbol{\Theta}(t)^{\top}, \bar{\boldsymbol{X}}_{2} \rangle \\
\vdots \\
\frac{d}{dt} \langle \boldsymbol{\Theta}(t) \boldsymbol{\Theta}(t)^{\top}, \bar{\boldsymbol{X}}_{N} \rangle
\end{bmatrix}$$

$$= 2 (F(\boldsymbol{\Theta}(t)) - \boldsymbol{y})^{\top} \begin{bmatrix}
\langle \bar{\boldsymbol{X}}_{1} \boldsymbol{\Theta}(t), \dot{\boldsymbol{\Theta}}(t) \rangle \\
\langle \bar{\boldsymbol{X}}_{2} \boldsymbol{\Theta}(t), \dot{\boldsymbol{\Theta}}(t) \rangle \\
\vdots \\
\langle \bar{\boldsymbol{X}}_{N} \boldsymbol{\Theta}(t), \dot{\boldsymbol{\Theta}}(t) \rangle^{\top}
\end{bmatrix}$$

$$= 2 (F(\boldsymbol{\Theta}(t)) - \boldsymbol{y})^{\top} \begin{bmatrix}
\operatorname{vec} (\bar{\boldsymbol{X}}_{1} \boldsymbol{\Theta}(t))^{\top} \\
\operatorname{vec} (\bar{\boldsymbol{X}}_{2} \boldsymbol{\Theta}(t))^{\top} \\
\vdots \\
\operatorname{vec} (\bar{\boldsymbol{X}}_{N} \boldsymbol{\Theta}(t))^{\top}
\end{bmatrix}$$

$$= (F(\boldsymbol{\Theta}(t)) - \boldsymbol{y})^{\top} J(\boldsymbol{\Theta}(t)) \operatorname{vec} (\dot{\boldsymbol{\Theta}}(t)). \tag{75}$$

Here, the Jacobian matrix $J(\mathbf{\Theta}(t))$ is defined as:

$$J(\boldsymbol{\Theta}(t)) \triangleq \frac{\partial F(\boldsymbol{\Theta}(t))}{\partial \text{vec}(\boldsymbol{\Theta}(t))} = \begin{bmatrix} \text{vec} \left(\nabla_{\boldsymbol{\Theta}} \langle \boldsymbol{\Theta}(t) \boldsymbol{\Theta}(t)^{\top}, \bar{\boldsymbol{X}}_{1} \rangle \right)^{\top} \\ \text{vec} \left(\nabla_{\boldsymbol{\Theta}} \langle \boldsymbol{\Theta}(t) \boldsymbol{\Theta}(t)^{\top}, \bar{\boldsymbol{X}}_{2} \rangle \right)^{\top} \\ \vdots \\ \text{vec} \left(\nabla_{\boldsymbol{\Theta}} \langle \boldsymbol{\Theta}(t) \boldsymbol{\Theta}(t)^{\top}, \bar{\boldsymbol{X}}_{N} \rangle \right)^{\top} \end{bmatrix} = 2 \begin{bmatrix} \text{vec} \left(\bar{\boldsymbol{X}}_{1} \boldsymbol{\Theta}(t) \right)^{\top} \\ \text{vec} \left(\bar{\boldsymbol{X}}_{2} \boldsymbol{\Theta}(t) \right)^{\top} \\ \vdots \\ \text{vec} \left(\bar{\boldsymbol{X}}_{N} \boldsymbol{\Theta}(t) \right)^{\top} \end{bmatrix} \in \mathbb{R}^{N \times 2d^{2}}.$$

$$(77)$$

With the notations defined above, we state the following theorem:

Theorem E.2. Let the combined weight matrix be

$$\boldsymbol{\Theta} \triangleq \begin{bmatrix} \boldsymbol{A} \\ \boldsymbol{B}^\top \end{bmatrix} \in \mathbb{R}^{2d \times d},$$

and consider the loss function $\tilde{\mathcal{L}}$ defined in (72). Denote

$$\sigma_{\min} \triangleq \sigma_{\min}(J(\mathbf{\Theta}(0))), \quad \sigma_{\max} \triangleq \sigma_{\max}(J(\mathbf{\Theta}(0))).$$

If the initialization satisfies:

$$\tilde{\mathcal{L}}(\mathbf{\Theta}(0)) \le \frac{\sigma_{\min}^6}{1152d\sigma_{\max}^2},$$

then for every $t \geq 0$ the following hold:

$$\tilde{\mathcal{L}}(\mathbf{\Theta}(t)) \leq \tilde{\mathcal{L}}(\mathbf{\Theta}(0)) \exp\left(-\frac{1}{2}\sigma_{\min}^2 t\right),$$
$$\|\mathbf{\Theta}(t) - \mathbf{\Theta}(0)\|_F \leq \frac{6\sqrt{2}\sigma_{\max}}{\sigma_{\min}^2} \sqrt{\tilde{\mathcal{L}}(\mathbf{\Theta}(0))}.$$

The above theorem tells us that, if the model is initialized with a sufficiently small loss, the model's loss will converge to zero quickly, and the parameters will not move significantly from the initialization. With the above theorem, we can state the following corollary:

Corollary E.3. Suppose A and B are initialized as balanced, i.e.:

$$\boldsymbol{A}(0)^{\top}\boldsymbol{A}(0) = \boldsymbol{B}(0)\boldsymbol{B}(0)^{\top}.$$

Under the conditions of Theorem E.2, *for every singular index* $i \in [d]$ *and all* $t \geq 0$:

$$\sigma_i(\boldsymbol{A}(t)) = \sigma_i(\boldsymbol{B}(t))$$
 and $|\sigma_i(\boldsymbol{A}(t)) - \sigma_i(\boldsymbol{A}(0))| \le \frac{\sigma_{\min}}{4\sqrt{2d}}$.

Consequently, the stable rank of A(t) remains bounded below by

$$\frac{\|\boldsymbol{A}(t)\|_F^2}{\|\boldsymbol{A}(t)\|_2^2} \ge \left(\frac{\|\boldsymbol{A}(0)\|_F - \frac{\sigma_{\min}}{4\sqrt{2d}}}{\|\boldsymbol{A}(0)\|_2 + \frac{\sigma_{\min}}{4\sqrt{2d}}}\right)^2.$$

E.3.1 PROOF OF THEOREM E.2

We begin the proof of the theorem by noting that the Jacobian $J(\cdot)$ is a Lipschitz function, as stated in the following lemma:

Lemma E.9. The Jacobian matrix $J(\mathbf{W})$, as defined in (77), is \sqrt{d} -Lipschitz. Specifically, for any matrices $\mathbf{W}, \mathbf{V} \in \mathbb{R}^{2d \times d}$, the following inequality holds:

$$||J(\boldsymbol{W}) - J(\boldsymbol{V})|| \le \sqrt{d} ||\operatorname{vec}(\boldsymbol{W}) - \operatorname{vec}(\boldsymbol{V})||.$$
(78)

Proof. Note that for each n-th observation,

$$J_n(\mathbf{\Theta}) = 2\text{vec}\left(\bar{\mathbf{X}}_n\mathbf{\Theta}\right)^{\top}$$

$$= \text{vec}\left(\begin{pmatrix} 0 & \mathbf{X}_n \\ \mathbf{X}_n^{\top} & 0 \end{pmatrix} \begin{pmatrix} \mathbf{A} \\ \mathbf{B}^{\top} \end{pmatrix}\right)^{\top}$$

$$= \text{vec}\left(\begin{pmatrix} \mathbf{X}_n\mathbf{B}^{\top} \\ \mathbf{X}_n^{\top}\mathbf{A} \end{pmatrix}\right)^{\top} \in \mathbb{R}^{2d^2}.$$

Let M_l denote the l-th row of a matrix M, and let M_{l-1} denote its l-th column. We have

$$\begin{aligned} \|J_n(\mathbf{\Theta})\|_F^2 &= \|\boldsymbol{X}_n^{\top} \boldsymbol{A}\|_F^2 + \|\boldsymbol{X}_n \boldsymbol{B}^{\top}\|_F^2 \\ &= \|\boldsymbol{e}_{j_n} \boldsymbol{e}_{i_n}^{\top} \boldsymbol{A}\|_F + \|\boldsymbol{e}_{i_n} \boldsymbol{e}_{j_n}^{\top} \boldsymbol{B}^{\top}\|_F \\ &= \|\boldsymbol{A}_{i_n}\|_2^2 + \|\boldsymbol{B}_{\cdot,j_n}\|_2^2. \end{aligned}$$

Now, suppose we observe all entries, i.e., $N = d^2$. Then for any fixed n, $i_n = i_m$ can be satisfied for all $m \in [d]$, meaning each element of A is observed d times. Similarly, each element of B is also observed d times.

Therefore, we can upper bound the Frobenius norm of the Jacobian matrix by the Frobenius norm of the Jacobian under full observation:

$$||J(\boldsymbol{\Theta})||_F^2 \le \sum_{n=1}^{d^2} (||\boldsymbol{X}_n^\top \boldsymbol{A}||_F^2 + ||\boldsymbol{X}_n \boldsymbol{B}^\top||_F^2)$$

= $d(||\boldsymbol{A}||_F^2 + ||\boldsymbol{B}||_F^2)$
= $d||\boldsymbol{\Theta}||_F^2$.

By upper-bounding the spectral norm of the difference between two Jacobian matrices and applying the inequality above, we obtain:

$$||J(\mathbf{W}) - J(\mathbf{V})||^2 = ||J(\mathbf{W} - \mathbf{V})||^2$$

 $\leq ||J(\mathbf{W} - \mathbf{V})||_F^2$
 $\leq d||\mathbf{W} - \mathbf{V}||_F^2$,

which concludes the proof.

Next, we borrow a lemma from Telgarsky (2021), which states that for a Lipschitz function J, if we consider a sufficiently small neighborhood around the initialization $\Theta(0)$, then the singular values of the Jacobian $J(\Theta)$ remain close to those at initialization:

Lemma E.10 (Lemma 8.3 in Telgarsky (2021)). *If we suppose* $\|\operatorname{vec}(\Theta) - \operatorname{vec}(\Theta(0))\| \le \frac{\sigma_{\min}}{2\sqrt{d}}$, we have the following:

$$\sigma_{\min}(J(\mathbf{\Theta})) \geq \frac{\sigma_{\min}}{2}, \quad \sigma_{\max}(J(\mathbf{\Theta})) \leq \frac{3\sigma_{\max}}{2},$$

where we denote $\sigma_{\min} \triangleq \sigma_{\min}(J(\mathbf{\Theta}(0)))$, and $\sigma_{\max} \triangleq \sigma_{\max}(J(\mathbf{\Theta}(0)))$.

For simplicity, we denote θ as the vectorized version of Θ , i.e., $\theta \triangleq \text{vec}(\Theta)$. We define the time step τ , which is the first time step when the trajectory of $\theta(t)$ touches the boundary:

$$\tau \triangleq \inf_{t \ge 0} \left\{ t \mid \|\boldsymbol{\theta}(t) - \boldsymbol{\theta}(0)\| \ge \frac{\sigma_{\min}}{2\sqrt{d}} \right\}.$$

We now demonstrate the convergence of the loss when $t \in [0, \tau]$ using the following lemma.

Lemma E.11. For all $t \in [0, \tau]$, the loss defined in (72) converges as follows:

$$\tilde{\mathcal{L}}(\boldsymbol{\Theta}(t)) \leq \tilde{\mathcal{L}}(\boldsymbol{\Theta}(0)) \exp\left(-\frac{1}{2}\sigma_{\min}^2 t\right),$$

where we define $\sigma_{\min} \triangleq \sigma_{\min}(J(\mathbf{\Theta}(0)))$.

Proof. Recall that the time derivative of the loss can be written as follows, according to (76):

$$\dot{\tilde{\mathcal{L}}}(\boldsymbol{\Theta}(t)) = -\left(F(\boldsymbol{\Theta}(t)) - \boldsymbol{y}\right)^{\top} J(\boldsymbol{\Theta}(t)) \dot{\boldsymbol{\theta}}(t)
= -\left(F(\boldsymbol{\Theta}(t)) - \boldsymbol{y}\right)^{\top} J(\boldsymbol{\Theta}(t)) J(\boldsymbol{\Theta}(t))^{\top} \left(F(\boldsymbol{\Theta}(t)) - \boldsymbol{y}\right),$$

noting that

$$\dot{\boldsymbol{\theta}}(t) = -\nabla_{\boldsymbol{\theta}(t)} \tilde{\mathcal{L}}(\boldsymbol{\Theta}(t)) = -J(\boldsymbol{\Theta}(t))^{\top} (F(\boldsymbol{\Theta}(t)) - \boldsymbol{y}).$$

By Lemma E.10, for any $t \in [0, \tau]$, we can upper bound the above term as follows:

$$\dot{\tilde{\mathcal{L}}}(\boldsymbol{\Theta}(t)) \leq -\lambda_{\min} \left(J(\boldsymbol{\Theta}(t)) J(\boldsymbol{\Theta}(t))^{\top} \right) \| F(\boldsymbol{\Theta}(t)) - \boldsymbol{y} \|^{2}
\leq -\frac{1}{2} \sigma_{\min}^{2} \tilde{\mathcal{L}}(\boldsymbol{\Theta}(t)).$$

Applying Grönwall's inequality gives:

$$\tilde{\mathcal{L}}(\boldsymbol{\Theta}(t)) \leq \tilde{\mathcal{L}}(\boldsymbol{\Theta}(0)) \exp\left(-\frac{1}{2}\sigma_{\min}^2 t\right) \quad \text{for } t \in [0,\tau].$$

The above lemma shows that the loss decays rapidly to zero if $\theta(t)$ stays within a small neighborhood around the initialization. We now show that if the loss converges quickly near initialization, then $\theta(t)$ does not move far from its initial value:

Lemma E.12. Let $\sigma_{\min} \triangleq \sigma_{\min}(J(\mathbf{\Theta}(0)))$ and $\sigma_{\max} \triangleq \sigma_{\max}(J(\mathbf{\Theta}(0)))$. For all $t \in [0, \tau]$, the distance between the weight vector at time t and the initial weight vector is bounded by:

$$\|\boldsymbol{\theta}(t) - \boldsymbol{\theta}(0)\| \le \frac{6\sqrt{2}\sigma_{\max}}{\sigma_{\min}^2} \sqrt{\tilde{\mathcal{L}}(\boldsymbol{\Theta}(0))}.$$

Proof. We start by evaluating the distance between $\theta(t)$ and $\theta(0)$ using Lemma E.10:

$$\|\boldsymbol{\theta}(t) - \boldsymbol{\theta}(0)\| = \left\| \int_0^t \dot{\boldsymbol{\theta}}(s) \, \mathrm{d}s \right\|$$

$$= \int_0^t \|J(\boldsymbol{\Theta}(s))^\top \left(F(\boldsymbol{\Theta}(s)) - \boldsymbol{y} \right) \| \, \mathrm{d}s$$

$$\leq \int_0^t \sigma_{\max} (J(\boldsymbol{\Theta}(s))) \|F(\boldsymbol{\Theta}(s)) - \boldsymbol{y}\| \, \mathrm{d}s$$

$$\leq \frac{3}{2} \sigma_{\max} \int_0^t \|F(\boldsymbol{\Theta}(s)) - \boldsymbol{y}\| \, \mathrm{d}s.$$

By Lemma E.11, we know that the objective function $\mathcal{\tilde{L}}(\Theta)$ satisfies:

$$||F(\boldsymbol{\Theta}(t)) - \boldsymbol{y}||^2 \le ||F(\boldsymbol{\Theta}(0)) - \boldsymbol{y}||^2 \exp\left(-\frac{1}{2}\sigma_{\min}^2 t\right).$$

Taking the square root of both sides, we obtain:

$$||F(\boldsymbol{\Theta}(t)) - \boldsymbol{y}|| \le ||F(\boldsymbol{\Theta}(0)) - \boldsymbol{y}|| \exp\left(-\frac{1}{4}\sigma_{\min}^2 t\right).$$

Substituting this into the previous inequality:

$$\|\boldsymbol{\theta}(t) - \boldsymbol{\theta}(0)\| \le \frac{3}{2}\sigma_{\max} \|F(\boldsymbol{\Theta}(0)) - \boldsymbol{y}\| \int_0^t \exp\left(-\frac{1}{4}\sigma_{\min}^2 s\right) \, \mathrm{d}s$$
$$\le \frac{6\sigma_{\max}}{\sigma_{\min}^2} \|F(\boldsymbol{\Theta}(0)) - \boldsymbol{y}\|,$$

where we used the fact that:

$$\int_0^t \exp(-Cs) \, \mathrm{d}s \le \frac{1}{C}, \quad \text{for } C > 0.$$

By combining Lemmas E.11 and E.12, we obtain the following results:

$$\tilde{\mathcal{L}}(\mathbf{\Theta}(t)) \le \tilde{\mathcal{L}}(\mathbf{\Theta}(0)) \exp\left(-\frac{1}{2}\sigma_{\min}^2 t\right),$$
 (79)

$$\|\boldsymbol{\theta}(t) - \boldsymbol{\theta}(0)\| \le \frac{6\sqrt{2}\sigma_{\max}}{\sigma_{\min}^2} \sqrt{\tilde{\mathcal{L}}(\boldsymbol{\Theta}(0))},$$
 (80)

which hold for $t \in [0, \tau]$. If we can demonstrate that $\tau = \infty$, the proof is complete.

Actually, if we initialize $\Theta(0)$ to satisfy the condition:

$$\tilde{\mathcal{L}}(\mathbf{\Theta}(0)) \le \frac{\sigma_{\min}^6}{1152d\sigma_{\max}^2},$$

and substitute this condition into (80), we obtain an upper bound for $\|\boldsymbol{\theta}(t) - \boldsymbol{\theta}(0)\|$:

$$\|\boldsymbol{\theta}(t) - \boldsymbol{\theta}(0)\| \le \frac{6\sqrt{2}\sigma_{\max}}{\sigma_{\min}^2} \frac{\sigma_{\min}^3}{\sqrt{1152d}\sigma_{\max}} = \frac{\sigma_{\min}}{4\sqrt{d}}.$$

Recall the definition of τ , which is the first time when $\theta(t)$ touches the boundary of the small ball around the initialization:

$$\tau \triangleq \inf_{t \ge 0} \left\{ t \mid \|\boldsymbol{\theta}(t) - \boldsymbol{\theta}(0)\| \ge \frac{\sigma_{\min}}{2\sqrt{d}} \right\}.$$

However, with the condition $\tilde{\mathcal{L}}(\Theta(0)) \leq \frac{\sigma_{\min}^6}{1152d\sigma_{\max}^2}$, $\boldsymbol{\theta}(t)$ cannot ever touch the boundary. This is because $\|\boldsymbol{\theta}(t) - \boldsymbol{\theta}(0)\|$ is bounded above by $\frac{\sigma_{\min}}{4\sqrt{d}}$, which is strictly less than $\frac{\sigma_{\min}}{2\sqrt{d}}$. Therefore, the parameter will remain inside the ball indefinitely, meaning $\tau = \infty$. This completes the proof of the theorem.

E.3.2 Proof of Corollary E.3

First, we establish the equality $\sigma_i(\boldsymbol{A}(t)) = \sigma_i(\boldsymbol{B}(t))$ for all $i \in [d]$. Corollary E.3 assumes that $\boldsymbol{A}(0)$ and $\boldsymbol{B}(0)$ are initialized as "balanced", satisfying $\boldsymbol{A}(0)^{\top}\boldsymbol{A}(0) = \boldsymbol{B}(0)\boldsymbol{B}(0)^{\top}$. By Lemma F.4, this balanced condition ensures that the singular values of $\boldsymbol{A}(t)$ and $\boldsymbol{B}(t)$ remain identical for all $t \geq 0$:

$$\sigma_i(\mathbf{A}(t)) = \sigma_i(\mathbf{B}(t)).$$

Second, we address the change in the singular values of a combined parameter matrix $\Theta(t)$ (related to A(t) and B(t)). Theorem E.2 states that under a specified condition on the initial loss, $\tilde{\mathcal{L}}(\Theta(0)) \leq \frac{\sigma_{\min}^6}{1152d\sigma_{\max}^2}$, the deviation of $\Theta(t)$ from its initialization $\Theta(0)$ is bounded for all $t \geq 0$ by:

$$\|\mathbf{\Theta}(t) - \mathbf{\Theta}(0)\|_F \le \frac{\sigma_{\min}}{4\sqrt{d}}.$$

Let $K = \frac{\sigma_{\min}}{4\sqrt{d}}$. By Weyl's inequality, $|\sigma_i(\boldsymbol{X}) - \sigma_i(\boldsymbol{Y})| \le ||\boldsymbol{X} - \boldsymbol{Y}||_2$, and noting that $||\cdot||_2 \le ||\cdot||_F$, we have for all $i \in [d]$:

$$|\sigma_i(\mathbf{\Theta}(t)) - \sigma_i(\mathbf{\Theta}(0))| \le ||\mathbf{\Theta}(t) - \mathbf{\Theta}(0)||_2$$

$$\le ||\mathbf{\Theta}(t) - \mathbf{\Theta}(0)||_F$$

$$\le K.$$

This inequality allows us to establish bounds for $\|\Theta(t)\|_F$ (using reverse triangle inequality) and its largest singular value $\sigma_1(\Theta(t)) = \|\Theta(t)\|_2$:

$$\|\mathbf{\Theta}(t)\|_F \ge \|\mathbf{\Theta}(0)\|_F - K,$$

$$\sigma_1(\mathbf{\Theta}(t)) \le \sigma_1(\mathbf{\Theta}(0)) + K.$$

This yields the following lower bound on the stable rank of $\Theta(t)$:

$$\frac{\|\mathbf{\Theta}(t)\|_F^2}{\|\mathbf{\Theta}(t)\|_2^2} \ge \left(\frac{\|\mathbf{\Theta}(0)\|_F - K}{\sigma_1(\mathbf{\Theta}(0)) + K}\right)^2 = \left(\frac{\|\mathbf{\Theta}(0)\|_F - \frac{\sigma_{\min}}{4\sqrt{d}}}{\|\mathbf{\Theta}(0)\|_2 + \frac{\sigma_{\min}}{4\sqrt{d}}}\right)^2.$$

Furthermore, the balancedness condition implies $\mathbf{A}(t)^{\top}\mathbf{A}(t) = \mathbf{B}(t)\mathbf{B}(t)^{\top}$. By the definition of $\mathbf{\Theta}(t)$, $\mathbf{\Theta}(t)^{\top}\mathbf{\Theta}(t) = \mathbf{A}(t)^{\top}\mathbf{A}(t) + \mathbf{B}(t)\mathbf{B}(t)^{\top}$, this leads to $\mathbf{\Theta}(t)^{\top}\mathbf{\Theta}(t) = 2\mathbf{A}(t)^{\top}\mathbf{A}(t)$. This relationship implies $\sigma_i(\mathbf{\Theta}(t)) = \sqrt{2}\sigma_i(\mathbf{A}(t))$ for all i. Substituting this into the bounds for $\mathbf{\Theta}(t)$, we have

$$\|\mathbf{A}(t)\|_F \ge \|\mathbf{A}(0)\|_F - K/\sqrt{2},$$

 $\|\mathbf{A}(t)\|_2 \le \|\mathbf{A}(0)\|_2 + K/\sqrt{2}.$

This leads to the final lower bound on the stable rank of A(t) (which, by balancedness, is equal to that of B(t)):

$$\frac{\|\boldsymbol{A}(t)\|_F^2}{\|\boldsymbol{A}(t)\|_2^2} \ge \left(\frac{\|\boldsymbol{A}(0)\|_F - K/\sqrt{2}}{\|\boldsymbol{A}(0)\|_2 + K/\sqrt{2}}\right)^2 = \left(\frac{\|\boldsymbol{A}(0)\|_F - \frac{\sigma_{\min}}{4\sqrt{2d}}}{\|\boldsymbol{A}(0)\|_2 + \frac{\sigma_{\min}}{4\sqrt{2d}}}\right)^2.$$

F USEFUL LEMMAS

Lemma F.1 (Adaptation of Lemma 1 and Theorem 3 in Arora et al. (2019)). For any time t, the product matrix $W(t) \in \mathbb{R}^{d,d}$ can be decomposed into its singular value decomposition:

$$oldsymbol{W}(t) = \sum_{r=1}^d \sigma_r(t) oldsymbol{u}_r(t) oldsymbol{v}_r(t)^ op$$

where $\sigma_r(t)$ are the singular values of W(t), and $u_r(t)$, $v_r(t)$ are the corresponding left and right singular vectors, respectively. Moreover, if A, B are balanced at initialization, i.e.,

$$\mathbf{A}^{\top}(0)\mathbf{A}(0) = \mathbf{B}(0)\mathbf{B}^{\top}(0),$$

the time evolution of the singular values $\sigma_r(t)$ is represented as:

$$\dot{\sigma_r}(t) = -2 \cdot \sigma_r(t) \cdot \left\langle \nabla \ell(\boldsymbol{W}(t)), \boldsymbol{u}_r(t) \boldsymbol{v}_r(t)^\top \right\rangle, \quad r = 1, \dots, d$$
 (81)

Lemma F.2. For any real-valued square matrix $A \in \mathbb{R}^{d \times d}$, the absolute value of its determinant equals the product of its singular values:

$$|\det(\mathbf{A})| = \prod_{r=1}^{d} \sigma_r$$

where σ_r are the singular values of **A**.

Proof. We express A using SVD: $A = U\Sigma V^{\top}$. Applying the determinant to both sides, we get:

$$\begin{aligned} \det(\boldsymbol{A}) &= \det(\boldsymbol{U}\boldsymbol{\Sigma}\boldsymbol{V}^{\top}) \\ &= \det(\boldsymbol{U})\det(\boldsymbol{\Sigma})\det(\boldsymbol{V}^{\top}) \end{aligned}$$

Here, U and V have orthonormal columns, and Σ is diagonal with singular values along its main diagonal. Since the determinant of an orthonormal matrix is either ± 1 ,

$$|\det(\boldsymbol{A})| = \det(\boldsymbol{\Sigma}) = \prod_{r=1}^d \sigma_r.$$

Lemma F.3 (Determinant of A(t)). Consider a matrix $A(t) \in \mathbb{R}^{d,d}$ initialized as $\det(A(0)) > 0$. Then, $\det(A(t)) > 0$ for all $t \geq 0$.

Proof. This follows directly from Lemma F.1 and F.2. Since the singular values are initialized as positive, and their evolution is continuous according to the given differential equation, they cannot become zero or negative. Therefore, A(t) maintains its sign of the determinant at initialization throughout the optimization process.

Lemma F.4 (Adaptation of Lemma 8 in Razin & Cohen (2020)). Consider a product matrix $W(t) = A(t)B(t) \in \mathbb{R}^{d \times d}$, where A(t) and B(t) are of equal size and balanced at initialization. Under these conditions, the following equality holds for all $t \geq 0$ and all singular values:

$$\sigma_r(\mathbf{W}(t)) = \sigma_r(\mathbf{A}(t))^2 = \sigma_r(\mathbf{B}(t))^2$$

where $\sigma_r(\cdot)$ denotes the r-th singular value of the respective matrix where $r \in [d]$. Moreover, if $\det(\mathbf{A}(0))$ and $\det(\mathbf{B}(0))$ are both positive, then by Lemma F.3, we can guarantee that for all $t \geq 0$:

$$\det \left(\boldsymbol{W}(t) \right) = \det \left(\boldsymbol{A}(t) \right)^2 = \det \left(\boldsymbol{B}(t) \right)^2$$

Lemma F.5 (Adaptation of Theorem 1 in Arora et al. (2019)). Consider a product matrix $W(t) = A(t)B(t) \in \mathbb{R}^{d \times d}$. We can guarantee A(t) and B(t) are analytic functions of t. As a result, W(t) is also an analytic function of t.

Lemma F.6 (Lemma 10 in Razin & Cohen (2020)). Let $f, g : [0, \infty] \to \mathbb{R}$ be real analytic functions such that $f^{(k)}(0) = g^k(0)$ for all $k \in \mathbb{N} \cup \{0\}$. Then, f(t) = g(t) for all $t \ge 0$.

Lemma F.7 (Positive Semidefiniteness of ABA^{\top}). For matrices $A, B \in \mathbb{R}^{d,d}$, if B is positive semi-definite, then both ABA^{\top} and $A^{\top}BA$ are positive semi-definite.

Proof. For any vector $\boldsymbol{x} \in \mathbb{R}^d$:

$$\boldsymbol{x}^{\top} \boldsymbol{A} \boldsymbol{B} \boldsymbol{A}^{\top} \boldsymbol{x} = (\boldsymbol{A}^{\top} \boldsymbol{x})^{\top} \boldsymbol{B} (\boldsymbol{A}^{\top} \boldsymbol{x}) \geq 0$$

since B is a positive semi-definite matrix. In the same way, for any vector $x \in \mathbb{R}^d$ we have:

$$\boldsymbol{x}^{\top} \boldsymbol{A}^{\top} \boldsymbol{B} \boldsymbol{A} \boldsymbol{x} = (\boldsymbol{A} \boldsymbol{x})^{\top} \boldsymbol{B} (\boldsymbol{A} \boldsymbol{x}) \geq 0$$

which concludes the proof.