

---

# Structured Identity Mapping: A Model for Out-of-Distribution Generalization Dynamics

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Modern multi-modal generative models exhibit remarkable out-of-distribution generalization capabilities, combining concepts in ways not observed in the training data. While a rich body of literature theoretically studies learning dynamics in in-distribution generalization settings, the dynamics of out-of-distribution generalization remain underexplored. In this work, we introduce and analyze the **Structured Identity Mapping** task, demonstrating how this simple model yields rich learning dynamics. Specifically, we analyze a one-hidden-layer network learning the identity map, using a training set composed of Gaussian point clouds structurally positioned at nodes of concept graphs. Our analysis of this model yields solutions that explain various empirical observations previously reported in text-conditioned diffusion models, including: (i) wave-like progression of compositional generalization dynamics, respecting hierarchical compositional structures; (ii) the impact of concept-centric data structures on concept learning speed; and (iii) non-monotonic progress of out-of-distribution generalization. In conclusion, our analytical model of concept learning establishes a theoretical foundation for investigating the dynamics of concept acquisition and combination in generative models.

## 17 1 Introduction

18 Concept learning and compositional generalization are essential features of modern generative models [1, 2, 3, 4, 5, 6]. These models can learn abstract concepts like shape, size and color from a limited set of training data and use them to generate images with novel combinations of concepts. In this paper, we focus on prompt conditioned image generative models, where the model is trained to generate images given a conditioning, e.g. a text prompt. Many papers have shown that this kind of model, especially ones based on diffusion models, often show some extent of concept learning [7, 8, 9].

24 One natural question which arises is: How do these capabilities emerge from training? It is generally hard to track the model behavior in terms of concept learning throughout training. Park et al. [6] proposed a principled way of studying it: instead of considering learning dynamics in parameter space as most previous work do, the authors study the learning dynamics in “concept space”. Briefly speaking, concept space is vector space that serves as an abstraction of real concepts. For each concept (e.g. color), a binary number can be used to represent its value (e.g. 0 for red and 1 for blue). In this way, a binary string can be mapped to a text condition (e.g. (1, 0, 1) might represent “big blue triangle”) and then be fed into the generative model as a conditioning vector. After that, a pre-trained classifier is used to check whether the model has indeed generated the specified combination of concepts. The idea of concept space is illustrated in Figure 1 as well as the middle figure in Figure 2.

34 Park et al. [6] found some interesting phenomena from the training dynamics in the concept space, such as control variables for concept learning and non monotonic trajectories. However, they only gave a description of these phenomena and thus the underlying causal mechanisms are still unclear. In this paper, we make a first step of establishing a theoretical framework that explains those phenomena.

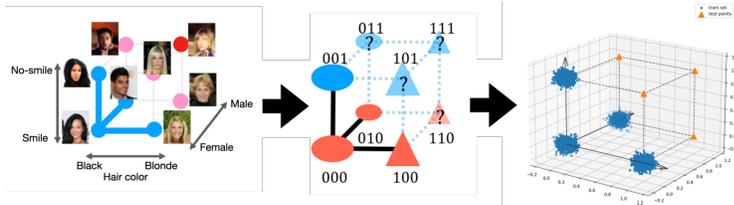
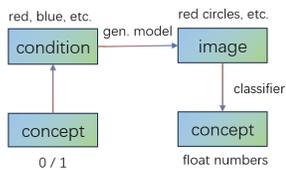


Figure 1: An schematic description of Concept Space. Figure 2: The abstraction process. Left: Actual Tasks (images from [10]) Middle: Concept Space considered in [6]; Right: SIM task considered in this work.

38 We would like to argue that the concept space framework established a very important idea: **In**  
 39 **concept space, generation is essentially identity mapping**, as the final classifier output should  
 40 match the input concept space specification. To this end, we make a further abstraction of the concept  
 41 learning process: learning identity mapping in Euclidean space. Each concept is represented as a  
 42 Gaussian cluster centered on an axes and the test point is a combination of several clusters means  
 43 (see the right figure in Figure 2). We call this task the Structured Identity Mapping learning (SIM)  
 44 task. As an analogue of concept learning, we claim that the seemingly simple SIM task has those key  
 45 features:

- 46 1. It is an out-of-distribution task, and theoretically non-generalizing solutions exists, so it is  
 47 non-trivial that a model solves this task;
- 48 2. As we will show later, this task captures many phenomena observed on real datasets;
- 49 3. This task masks out distracting terms and is simple enough that we can elaborate the  
 50 behaviour of the model on it.

51 In this work, we give a detailed description of how the model behaves throughout training on the  
 52 SIM task and how those behaviors corresponds to the behavior of actual diffusion models. After that,  
 53 we study the learning dynamics of a specific model on this task, and give a rigorous proof of the  
 54 existence of the phenomena. More specifically, we make two major contributions:

- 55 1. In Section 3, we train MLP regression models on the SIM task and empirically reveal some  
 56 key features of the training dynamics, including: 1) how the order of concepts learning  
 57 is controlled by the signal strength and diversity of the dataset; 2) the deceleration of  
 58 concept learning with training progress and 3) a “Transient Memorization” phenomena  
 59 where the model shows a trend of generalization at the beginning of the training soon  
 60 followed by memorization before heading back to generalization again. The last phenomena  
 61 leads to a double descent-like loss curve with respect to optimization steps even in an  
 62 under-parameterized setting.
- 63 2. In Appendix A, we explain the empirical phenomena by a detailed analysis of two simplified  
 64 models. Interestingly, we show that *Transient Memorization* is a phenomenon only observed  
 65 in multi-layer models, which reveals a key difference between one-layer model and deep  
 66 models. Our novel analysis of dynamics reveals a multi-stage evolution process of the  
 67 Jacobian of a two layer symmetric linear model ( $f(x; U) = UU^T x$ ), and we show that  
 68 each stage of the Jacobian evolution precisely corresponds to the stages of the *Transient*  
 69 *Memorization*.

## 70 2 Preliminaries and Problem Setting

71 Throughout the paper, we use bold lowercase letters (e.g.  $\mathbf{x}$ ) to represent vectors, and bold uppercase  
 72 letters (e.g.  $\mathbf{A}$ ) for matrices. Non-bold versions with subscripts represent corresponding entries of the  
 73 vectors or matrices, e.g.  $x_i$  represent the  $i$ -th entry of  $\mathbf{x}$  and  $a_{i,j}$  represent the  $(i, j)$ -th entry of  $\mathbf{A}$ .

74  $[a]$  represents the set of all natural numbers that is smaller or equal to  $a$ , i.e.  $[a] = \{1, 2, \dots, a\}$ .

75  $\mathbf{1}_k$  represents a one-hot vector which is 0 at every entry except the  $k$ -th entry being 1. The dimen-  
 76 sionality of the vector is determined by the context if not specified.  $\mathbf{I}_k$  represents a diagonal matrix  
 77 whose first  $k$  diagonal entries are 1 and others are 0.

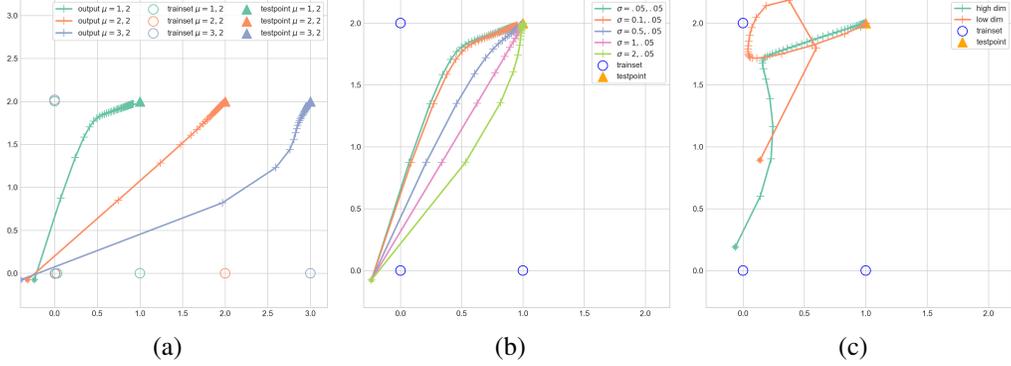


Figure 3: The two dimensional ( $s = 2$ ) output dynamics under different settings, evaluated for the test point,  $(1, 1)$ . We only show the center of the training classes as a circle, but the actual training set can have varied shapes based on the configuration of  $\sigma$ . (a) one layer linear model with  $\sigma_{:2} = (.05, .05)$  and varying  $\mu$ ; (b) one layer linear model with  $\mu_{:2} = (1, 2)$  and varying  $\sigma$ ; (c) 4 layer linear models for different model dimensions. high dim:  $d = 64$ , low dim:  $d = 2$ . Note that (a) and (b) are both in high dimensional model setting.

## 78 2.1 Problem Setting

79 **Data.** The SIM dataset is composed of several Gaussian clusters, each occupying a coordinate  
80 direction. Figure 2 (right) illustrates the SIM dataset.

81 Let  $d \in \mathbb{N}$  be the dimensionality of the input space,  $s \in [d]$  be the number of clusters, and  $n \in \mathbb{N}$  be  
82 the number of samples from each cluster. The training set  $\mathcal{D} = \bigcup_{p \in [s]} \{ \mathbf{x}_k^{(p)} \}_{k=1}^n$  is generated by  
83 the following process: for each  $p \in [s]$ , the data sample is sampled i.i.d. from a Gaussian distribution  
84  $x_k^{(p)} \sim \mathcal{N} \left[ \mu_p \mathbf{1}_p, \text{diag}(\boldsymbol{\sigma})^2 \right]$ , where  $\mu_p \geq 0$  is the distance of the  $p$ -th cluster center from the origin,  
85 and  $\boldsymbol{\sigma}$  is a vector with only the first  $s$  entries being non-zero, where  $\sigma_i$  represents the variance on the  
86  $i$ -th direction. Notice that we allow  $\mu_p = 0$  for a specific  $p$  to create a cluster that centers at  $\mathbf{0}$  while  
87 keeping the terminology as simple as possible.

88 **Loss function.** The training problem is to learn identity mapping on  $\mathbb{R}^d$ . For a model  $f : \mathbb{R}^m \times \mathbb{R}^d \rightarrow$   
89  $\mathbb{R}^d$  and a parameter vector  $\boldsymbol{\theta} \in \mathbb{R}^m$ , we train the model parameter  $\boldsymbol{\theta}$  with the mean square error loss.

$$\mathcal{L}(\boldsymbol{\theta}) = \frac{1}{2sn} \sum_{p=1}^s \sum_{k=1}^n \left\| f(\boldsymbol{\theta}; \mathbf{x}_k^{(s)}) - \mathbf{x}_k^{(s)} \right\|^2. \quad (2.1)$$

90 **Evaluation.** In the main paper, we focus on a single test point  $\hat{\mathbf{x}} = \sum_{p=1}^s \mu_p \mathbf{1}_p$ . Notice that this  
91 point is outside of the training distribution. In Appendix B, we report additional results for the case  
92 of multiple test points.

## 93 3 Observations on the SIM Task

94 In this section, we summarize some of the interesting empirical findings of the model behavior  
95 when we conduct experiments on the SIM. We note that these findings reported in this section are in  
96 one-to-one correspondence with results with diffusion models.

97 In all experiments, we use MLP models. We perform experiments with both linear activations and  
98 non-linear ReLU activations, as well as model with different number of layers. Due to the limited  
99 space, we only show the results of some settings here but the phenomena reported are consistent with  
100 different hyperparameters.

### 101 3.1 Generalization Order is Controlled by Signal Strength and Diversity

102 One interesting findings from previous work is that if we alter the strength of one signal from small  
103 to large, the concur of the learning dynamics would dramatically change [6]. Moreover, it is also  
104 commonly hypothesised that with more diverse data, the model should also generalize better [11, 12].

105 In the SIM task the distance  $\mu_k$  of each cluster can be viewed as the corresponding signal strength,  
 106 and the covariance  $\sigma_k$  can be viewed as the data diversity. In Figure 3, we train models with  $s = 2$   
 107 allowing us to directly plot the trajectory of the model output at each timestep. In this case, there are  
 108 two components,  $x$  and  $y$ , to be learned and the order of learning can be seen from the trajectory.

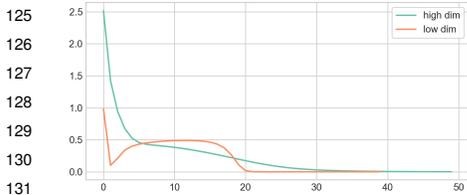
109 In Figure 3 (a)<sup>1</sup>,  $\sigma$  is fixed, and we can see when  $\mu_1 < \mu_2$ , the dynamics shows an upward bulging,  
 110 showing a preference for the direction of stronger signal. As we gradually increase  $\mu_1$ , this trajectory  
 111 gradually transitions from an upward bulge to downward one, consistent with the stronger signal  
 112 strength.

113 In Figure 3 (b), the  $\mu$  is fixed to have one signal stronger than the other, and the model, as expected,  
 114 prefers the direction with stronger signal when the data diversity is balanced. However, if we tune  
 115 the data diversity of one side from weak to strong while keeping the other side unchanged, we can  
 116 override the preference coming from the mean signal.

117 The results in Figure 3 (a) and (b) gives us a very concrete conclusion: The learning direction is a  
 118 competition driven by signal strength and diversity, and the model prefers direction that has stronger  
 119 signal and more diversity.

### 120 3.2 Transient Memorization

121 The results in Figure 3 (a) and (b) are both performed with one layer models and under a high  
 122 dimensional setting ( $d = 64$ ). Despite the overall trend is similar in other settings, it is worth  
 123 exploring the change of trajectory as we increase the number of layers, and / or reduce the dimension.  
 124



133 Figure 4: The loss function of multi layer  
 134 models.

In Figure 3 (c), we perform experiments with deeper models, and optionally with a lower dimension. Under these changes, we find that the model shows an interesting irregular behavior, where it initially heads towards the right direction, but soon turns toward the training set cluster with the strongest signal, exhibiting distributional memorization the training set. However, with enough training, the model correctly moves towards the intended target and thus generalizes. We call this behavior **Transient Memorization**.

135 This trajectory could be suggestive of a non-monotonic  
 136 generalization. We track the value of the loss function during training in Figure 4, demonstrating  
 137 a double descent-like curve. We note that the *Transient Memorization* phenomena seems to be  
 138 strongest when the dimensionality is low, and is rather modest with high dimensional settings. In  
 139 the high dimensional setting the loss descent slows down at some point but doesn't actually exhibit  
 140 non-monotonic behavior.. This low dimensional preference can also be explained perfectly by our  
 141 theory, further described in Appendix A.

### 142 3.3 Convergence Rate Slow Down In Terminal Phase

143 In Figure 3, the markers on the curve corresponds to equal time intervals. One can observe that the  
 144 model's generation's evolution in concept space slows down as training progresses. This phenomena  
 145 is observed in diffusion models as well [6].

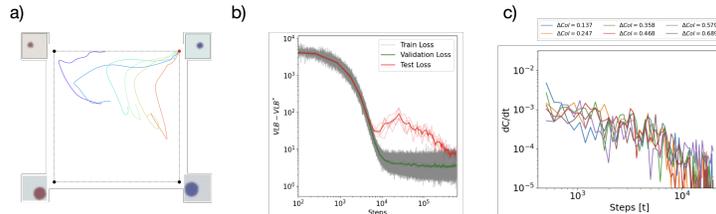


Figure 5: **Diffusion Model Results** a) Concept signal controls learning order and speed b) Transient Memorization observed in vector space diffusion models c) Deceleration of concept learning with time

<sup>1</sup>In Figure 3 (a) some training set centroids are overlapping and we perturb them a little for visibility.

## References

- 146
- 147 [1] Maya Okawa, Ekdeep S Lubana, Robert Dick, and Hidenori Tanaka. Compositional abilities  
148 emerge multiplicatively: Exploring diffusion models on a synthetic task. *Advances in Neural*  
149 *Information Processing Systems*, 36, 2024.
- 150 [2] Parikshit Ram, Tim Klinger, and Alexander G Gray. What makes models compositional? a  
151 theoretical view: With supplement. *arXiv preprint arXiv:2405.02350*, 2024.
- 152 [3] Thaddäus Wiedemer, Prasanna Mayilvahanan, Matthias Bethge, and Wieland Brendel. Com-  
153 positional generalization from first principles. *Advances in Neural Information Processing*  
154 *Systems*, 36, 2024.
- 155 [4] Xindi Wu, Dingli Yu, Yangsibo Huang, Olga Russakovsky, and Sanjeev Arora. Conceptmix:  
156 A compositional image generation benchmark with controllable difficulty, 2024. URL <https://arxiv.org/abs/2408.14339>.  
157
- 158 [5] Nan Liu, Shuang Li, Yilun Du, Antonio Torralba, and Joshua B. Tenenbaum. Compositional  
159 visual generation with composable diffusion models, 2023. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2206.01714)  
160 [2206.01714](https://arxiv.org/abs/2206.01714).
- 161 [6] Core Francisco Park, Maya Okawa, Andrew Lee, Ekdeep Singh Lubana, and Hidenori Tanaka.  
162 Emergence of hidden capabilities: Exploring learning dynamics in concept space. *arXiv preprint*  
163 *arXiv:2406.19370*, 2024.
- 164 [7] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark  
165 Chen, and Ilya Sutskever. Zero-shot text-to-image generation, 2021.
- 166 [8] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer.  
167 High-resolution image synthesis with latent diffusion models, 2022.
- 168 [9] Core Francisco Park, Maya Okawa, Andrew Lee, Ekdeep Singh Lubana, and Hidenori Tanaka.  
169 Emergence of hidden capabilities: Exploring learning dynamics in concept space, 2024. URL  
170 <https://arxiv.org/abs/2406.19370>.
- 171 [10] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Large-scale celebfaces attributes  
172 (celeba) dataset. *Retrieved August*, 15(2018):11, 2018.
- 173 [11] Zhiqiang Gong, Ping Zhong, and Weidong Hu. Diversity in machine learning. *Ieee Access*, 7:  
174 64323–64350, 2019.
- 175 [12] José F Díez-Pastor, Juan J Rodríguez, César I García-Osorio, and Ludmila I Kuncheva. Diversity  
176 techniques improve the performance of the best imbalance learning ensembles. *Information*  
177 *Sciences*, 325:98–117, 2015.
- 178 [13] Sanjeev Arora, Nadav Cohen, Wei Hu, and Yuping Luo. Implicit regularization in deep matrix  
179 factorization. *Advances in Neural Information Processing Systems*, 32, 2019.
- 180 [14] Sanjeev Arora, Nadav Cohen, Noah Golowich, and Wei Hu. A convergence analysis of gradient  
181 descent for deep linear neural networks. *arXiv preprint arXiv:1810.02281*, 2018.
- 182 [15] Ziwei Ji and Matus Telgarsky. Gradient descent aligns the layers of deep linear networks. *arXiv*  
183 *preprint arXiv:1810.02032*, 2018.
- 184 [16] Dominik Stöger and Mahdi Soltanolkotabi. Small random initialization is akin to spectral  
185 learning: Optimization and generalization guarantees for overparameterized low-rank matrix  
186 reconstruction. *Advances in Neural Information Processing Systems*, 34:23831–23843, 2021.
- 187 [17] Jikai Jin, Zhiyuan Li, Kaifeng Lyu, Simon Shaolei Du, and Jason D Lee. Understanding incre-  
188 mental learning of gradient descent: A fine-grained analysis of matrix sensing. In *International*  
189 *Conference on Machine Learning*, pages 15200–15238. PMLR, 2023.

## 190 A Theoretical Explanation

191 In this section, we study the training dynamics of a specific type of linear models which is tractable  
 192 on the SIM task, and we provide explanations of the behaviors on the SIM task. We first consider a  
 193 linear model and show that despite it can explain some of the phenomena, the linear model can not  
 194 actually reproduce every phenomenon, which suggests that some phenomena are intrinsic for deep  
 195 models, which highlights the difference of shallow and deep models.

196 Throughout this section, we assume  $f(\boldsymbol{\theta}; \mathbf{x})$  is a linear operator of  $\mathbf{x}$ . In this case the Jacobian of  $f$   
 197 w.r.t.  $\mathbf{x}$  is a matrix that is completely determined by  $\boldsymbol{\theta}$ , which we denote by  $\mathbf{W}_\theta = \frac{\partial f(\boldsymbol{\theta}; \mathbf{x})}{\partial \mathbf{x}}$ . It's easy  
 198 to see that we have  $f(\boldsymbol{\theta}; \mathbf{x}) = \mathbf{W}_\theta \mathbf{x}$ . Through the trace trick, it's easy to show that the overall loss  
 199 function is equal to

$$\mathcal{L}(\boldsymbol{\theta}) = \frac{1}{2} \left\| (\mathbf{W}_\theta - \mathbf{I}) \mathbf{A}^{1/2} \right\|_{\mathcal{F}}^2, \quad (\text{A.1})$$

200 where  $\mathbf{A} = \frac{1}{sn} \sum_{p=1}^s \sum_{k=1}^n \mathbf{x}_k^{(p)} \mathbf{x}_k^{(p)\top}$  is the empirical covariance. When  $n$  is large,  $\mathbf{A}$  converges  
 201 to the true covariance of the dataset  $\mathbf{A} \rightarrow \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \mathbf{x} \mathbf{x}^\top$ . To avoid extra distractions, through out this  
 202 section we assume  $\mathbf{A}$  equals the the true covariance, which is a diagonal matrix  $\mathbf{A} = \text{diag}(\mathbf{a})$  defined

203 by  $a_p = \begin{cases} \sigma_p^2 + \frac{\mu_p^2}{s} & p \leq s \\ 0 & p > s \end{cases}$ , for any  $p \in [s]$ . For completeness, we write the full derivation of  
 204 eq. (A.1) and  $\mathbf{A}$  in Appendix C.1.

205 **Remark.** Notice that in the linear setting we might not directly train  $\mathbf{W}_\theta$ , instead we train its  
 206 components. For example we might have  $\boldsymbol{\theta} = (\mathbf{W}_1, \mathbf{W}_2)$  and have  $\mathbf{W}_\theta = \mathbf{W}_1 \mathbf{W}_2$ . What we actually  
 207 train is  $\mathbf{W}_1$  and  $\mathbf{W}_2$  instead of  $\mathbf{W}_\theta$ . As many previous work have emphasized [13, 14, 15], although  
 208 the deep linear model has the same capacity as a one-layer linear model, their dynamics can be vastly  
 209 different and the loss landscape of deep linear models can be non-convex.

### 210 A.1 The Failure of One Layer Model Theory

211 As a warm-up, we first study the dynamics of one layer linear models, i.e.  $f(\mathbf{W}; \mathbf{x}) = \mathbf{W} \mathbf{x}$ , in which  
 212 case the Jacobian  $\mathbf{W}_\theta$  is simply  $\mathbf{W}$ . We will show that it can explain some of the phenomenons but  
 213 fail to capture other interesting behaviours.

214 **Theorem A.1.** Let  $\mathbf{W}(t) \in \mathbb{R}^{d \times d}$  be initialized as  $\mathbf{W}(0) = \mathbf{W}^{(0)}$ , and updated by

$$\frac{d\mathbf{W}(t)}{dt} = -\nabla \mathcal{L}(\mathbf{W}), \quad (\text{A.2})$$

215 with  $\mathcal{L}$  be defined by eq. (A.1) with  $f(\mathbf{W}, \mathbf{z}) = \mathbf{W} \mathbf{z}$ , then we have for any  $\mathbf{z} \in \mathbb{R}^d$ ,

$$f(\mathbf{W}(t), \mathbf{z})_k = \underbrace{\mathbb{1}_{\{k \leq s\}} [1 - \exp(-a_k t)] z_k}_{\tilde{G}_k(t)} + \underbrace{\sum_{i=1}^s \exp(-a_i t) w_{k,i}(0) z_i}_{\tilde{N}_k(t)}. \quad (\text{A.3})$$

216 See Appendix C.2 for a proof of Theorem A.1.

217 The output of the one-layer model  $f(\boldsymbol{\theta}(t); \mathbf{z})$  can be decomposed into two terms: the *growth term*  
 218  $\tilde{G}_k(t; \mathbf{z}) = \mathbb{1}_{\{k \leq s\}} [1 - \exp(-a_k t)] z_k$  and the *noise term*  $\tilde{N}_k(t; \mathbf{z}) = \sum_{i=1}^s \exp(-a_i t) w_{k,i}(0) z_i$ .  
 219 By observing these two terms we can find the following properties: 1) the growth term converges to  
 220  $x_k$  when  $k \leq s$  and 0 when  $k > s$ , and the noise term always converges to 0; 2) both terms converges  
 221 in an exponential rate; 3) the noise term is upper bounded by  $\sum_{i=1}^s w_{k,i}(0) x_i$ .

222 If the model is initialized small, specifically  $w_{k,i}(0) \ll \frac{1}{s \max_{i \in [s]} \{x_i\}}$ , then the  $\tilde{N}_k(t)$  will always be  
 223 small, and thus can be omitted. With this assumption in effect, the model output is dominated by the  
 224 growth term. A closer look at the growth term reveals the origin of some of the phenomena observed  
 225 before.

226 **Generalization Order.** In eq. (A.3) we can see that the growth term  $\tilde{G}_k(t; \mathbf{z})$  converges in an  
 227 exponential speed, with the exponential term controlled by  $a_k$ , which is  $a_k = \frac{s\sigma_k^2 + \mu_k^2}{s}$ . Therefore, the  
 228 direction with larger  $a_k$ , i.e. larger  $\mu_k$  and / or  $\sigma_k$ , converges faster. The equation also reveals the  
 229 proportional relationship of  $\mu_k$  and  $\sigma_k$ .

230 **Terminal Phase Slowing Down.** By taking derivative of the growth term  $\tilde{G}_k(t; \mathbf{z})$  with  $k \leq s$ ,  
 231 we have

$$\dot{\tilde{G}}_k(t; \mathbf{z}) = a_k z_k \exp(-a_k t), \quad (\text{A.4})$$

232 which monotonically decays w.r.t.  $t$ , and reveals an (exponentially) slowing down of the convergence  
 233 rate when  $t$  becomes large.

234 **The Failure of One Layer Model Theory.** We have shown that Theorem A.1 can explain many  
 235 observations. However, in eq. (A.3) the growth term  $\tilde{G}_k(t; \mathbf{z})$  is independent and monotonic for each  
 236 layer, which only produces monotonic and rather regular traces (this is verified by the experiments in  
 237 Section 3.1. However, as the experiments in Section 3.2 show, when the number of layer becomes  
 238 larger, the model actually shows a non-monotonic trace that can have detours. The theory based on  
 239 one layer model fails in capturing that phenomenon. In the subsequent section, we introduce a more  
 240 complex theory by considering a deeper model, and we will show that this theory explains all the  
 241 phenomena observed in Section 3.

## 242 A.2 A Symmetric Two Layer Model Theory

243 In this subsection, we introduce a two layer model with symmetric weights. Despite its simplicity,  
 244 we show that it perfectly captures every observations presented in Section 3. More importantly, the  
 245 theory derived from this model draws a clear picture of the evolution of the evolution of the model  
 246 Jacobian and provides us with a clear and understandable explanation of the origin of the seemingly  
 247 irregular behaviours of the model.

248 Due to the space limitation, in this subsection we focus on providing an intuitive explanation of the  
 249 model behaviour and delay the formal proof to later sections.

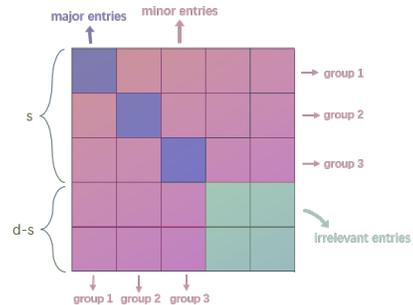
250 Formally, in this subsection, we consider a linear model that has two layers with shared weights,  
 251 namely

$$f(\mathbf{U}; \mathbf{x}) = \mathbf{U}\mathbf{U}^\top \mathbf{x}, \quad (\text{A.5})$$

252 where  $\mathbf{U} \in \mathbb{R}^{d \times d'}$  and  $d' > d$ . Notice that this is a model commonly studied in theory [16, 17], and  
 253 our analysis goes beyond the existing ones by providing an analysis for the early stage phenomenon,  
 254 Transient Memorization. For simplicity, in this subsection we denote the Jacobian of  $f$  at time point  $t$   
 255 by  $\mathbf{W}(t) = \mathbf{W}_{\mathbf{U}(t)}$ , then the update of the  $i, j$ -th entry of  $\mathbf{W}$  is given by

$$\dot{w}_{i,j} = \underbrace{w_{i,j}(a_i + a_j)}_{G_{i,j}(t)} - \underbrace{\frac{1}{2} w_{i,j} [w_{i,i}(3a_i + a_j) + \mathbb{1}_{\{i \neq j\}} w_{j,j}(3a_j + a_i)]}_{S_{i,j}(t)} - \underbrace{\frac{1}{2} \sum_{\substack{k \neq i \\ k \neq j}} w_{k,i} w_{k,j} (a_i + a_j + 2a_k)}_{N_{i,j}(t)}. \quad (\text{A.6})$$

256 As noted in eq. (A.6), we decompose the  
 257 update of  $w_{i,j}$  into three terms. We call  
 258  $G_{i,j}(t) = w_{i,j}(t)(a_i + a_j)$  the *growth term*,  $S_{i,j}(t) =$   
 259  $\frac{1}{2} w_{i,j} [w_{i,i}(3a_i + a_j) + \mathbb{1}_{\{i \neq j\}} w_{j,j}(3a_j + a_i)]$  the *sup-*  
 260 *pression term*, and  $N_{i,j}(t) = \frac{1}{2} \sum_{\substack{k \neq i \\ k \neq j}} w_{k,i}(t) w_{k,j}(t) (a_i +$   
 261  $a_j + 2a_k)$  the *noise term*. The name of the three terms  
 262 suggests their role in the evolution of the Jacobian: the  
 263 growth term  $G_{i,j}$  always has the same sign as  $w_{i,j}$ , and  
 264 has a positive contribution to the update, so it always leads  
 265 to the direction that **enlarges the absolute value** of  $w_{i,j}$ ;  
 266 the suppression term  $S_{i,j}$  also has the same sign<sup>2</sup> as  $w_{i,j}$ ,



<sup>2</sup>Notice that since  $\mathbf{W} = \mathbf{U}\mathbf{U}^\top$  is a PSD matrix, the diagonal entries are always non-negative.

Figure 6: An illustration of the entries of the Jacobian.

267 but has a negative contribution in the update function of  
 268  $w_{i,j}$ , so it always leads to the direction that **shrinks the**  
 269 **absolute value** of  $w_{i,j}$ ; the effect of noise term is random  
 270 since it depends on the sign of  $w_{i,j}$  and other terms. It is  
 271 proved in Lemma D.8 that under specific assumptions, the  
 272 noise term will never be too large so for the sake of brevity, we ignore it in the following discussion  
 273 and delay the treatment of it to the rigorous proofs in later sections.

### 274 A.2.1 The Evolution of Entries of Jacobian

275 In order to better present the evolution of the Jacobian, we divide the entries of the Jacobian into three  
 276 types: the **major entries** are the first  $s$  diagonal entries, and the **minor entries** are the off-diagonal  
 277 entries who are in the first  $s$  rows or first  $s$  columns, and other entries are **irrelevant entries**. Notice  
 278 that the irrelevant entries doesn't contribute to the output of the test points so we ignore them.  
 279 Moreover, we also divide minor entries to groups. The minor entry in the  $p$ -th row or column belongs  
 280 to the  $p$ -th group (each entry can belong to at most two groups). See Figure 6 for an illustration of the  
 281 division of the entries.

282 **Initial Growth.** In this section, we assume  $w_{i,i}$ -s are initialized around a very small value  $\omega$  such  
 283 that  $\omega \ll \frac{1}{d \max_{i \in [s]} a_i}$  (See Appendix D.1 for specific assumptions). One can easily notice that when  
 284 all  $w_{i,j}$ -s are around  $\omega$ , the growth term is  $o(\omega)$  and the suppression term and the noise term are both  
 285  $o(\omega^2)$ . This indicates when all entries are closed to initialization, the suppression term is negligible  
 286 and the evolution of  $w_{i,j}$  is dominated by the growth term. Therefore, in this stage, every value in the  
 287 Jacobian grows towards the direction of enlarging its absolute value, with the speed determined by  
 288  $a_i + a_j$ . Since we assumed that  $\mathbf{a}$  is ordered in a descending order, it's not hard to see that each entry  
 289 grows faster than all entries below or on the right of it. The initial growth stage is characterized by  
 290 Lemmas D.1 to D.3.

291 **First Suppression.** We say an entry that is close to its initialization is in the "initial phase". In the  
 292 Initial Growth stage, the first major entry will be the one that grows exponentially faster than all other  
 293 entries, thus it will be the first one that leaves the initial phase. Once the first major entry becomes  
 294 significant and non-negligible, it will effect on the suppression term of all minor entries in the first  
 295 group. When the difference between  $a_1$  and  $a_2$  is large enough, the first major entry is able to change  
 296 the growth direction of the first group of minor entries and push their value to 0. The suppression  
 297 stages are characterized by Lemma D.7

298 **Second Growth and Cycle.** Once the suppression of the first group of minor entries takes into  
 299 effect, the second major entry becomes the one that grows fastest. Thus, the second major entry will  
 300 be the second one that leaves the initial stage. When the second major entry becomes large enough,  
 301 again it will suppress the second group of minor entries and push their value to 0. The process  
 302 continues like this: A major entry growth is followed by the suppression of the corresponding group  
 303 of minor entries, and the suppression leaves space for the growth of the next major entry. The general  
 304 growth stages are characterized by Lemma D.4 and the fate of off-diagonal entries are characterized  
 305 by Lemma D.8.

306 **Growth Slow Down and Stop.** Notice that the suppression term of a major entry is also determined  
 307 by itself. Thus when a major entries becomes significantly large, it also suppresses itself, and causes  
 308 the slow down of growth. Note that this won't reverse its growth direction since the suppression term  
 309 is always smaller than the growth term, until  $w_{i,i}$  becomes one where the growth and suppression  
 310 equal and the evolution stops. The terminal stage of the growth of major entries are characterized by  
 311 Lemma D.5.

### 312 A.2.2 Explaining Model Behavior

313 Recall that we have  $f(U(t); \hat{\mathbf{x}})_k = \sum_{p=1}^s w_{k,p}(t) v_p \mu_p$ . In this subsection we explain how the  
 314 Jacobian evolution predicted in Appendix A.2.1 are reflected in the model output evolution.

315 **Learning Order and Terminal Slowing Down.** From the discussions in Appendix A.2.1, we see  
 316 that at the end of the learning, all the major entries converges to 1 and all minor entries converges to

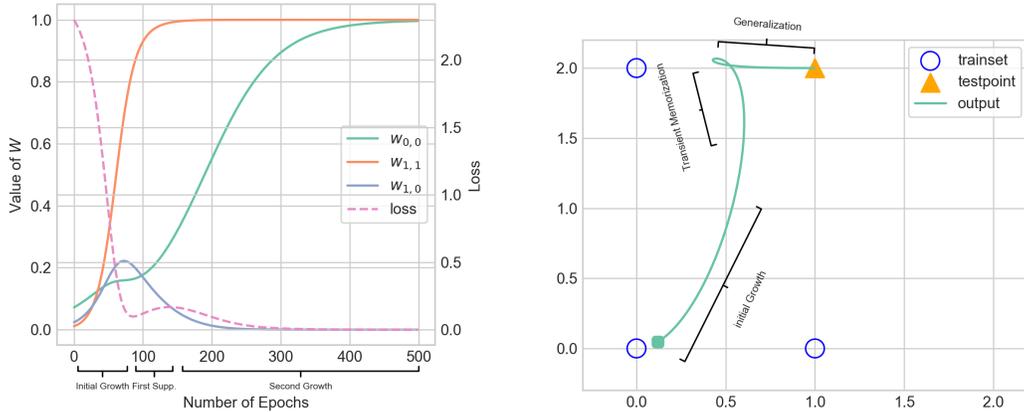


Figure 7: The learning dynamics of a two layer symmetric linear mode. Left: The change of the loss and the Jacobian entries with time predicted by the theory; Right: the corresponding model output curve. The figures are plotted under  $s = 2$  and all entries of  $\mathbf{W}$  are initialized positive.

317 0, and the major entries grows in the order of corresponding  $a_p$ , which depends on the  $\mu_p$  and  $\sigma_p$ ,  
 318 and slows down when approaching the terminal. This explains our observation that directions with  
 319 larger  $\mu_p$  and / or  $\sigma_p$  is learned first, as well as the terminal phase slowing down of the learning.

320 **Transient Memorization.** We argue that, the Transient Memorization is caused by the initial  
 321 growth. Notice that in the initial stage  $f(\mathbf{U}(t); \hat{\mathbf{x}})_k$  is dominated by  $w_{k,1}(t)v_1\mu_1$ , since  $w_{k,1}$  grows  
 322 fastest among all the entries. If  $w_{k,1}$  happens to be initialized positive, it growth towards the positive  
 323 direction. When  $s$  is small, this is actually easy to satisfy<sup>3</sup>. This causes an illusion that the model is  
 324 going towards the right direction, because the target point has all positive coordinates.

325 Figure 7 shows the loss curve and Jacobian entry evolution predicted by the theory with all entries of  
 326  $\mathbf{W}$  initialized positive. Notice that how the first and second descending of loss accurately corresponds  
 327 to the initial and second growth of the major entries, and the ascending of the loss corresponds to the  
 328 suppression of the minor entries.

329 **Remark.** We note that, for a major entry  $w_{i,i}$  in the Jacobian  $\mathbf{W}$ , Lemma D.4 proves that when  
 330  $w_{i,i} = \lambda$ , the growth rate of  $w_{i,i}$  is  $\lambda a_i \exp(-2\lambda a_k t)$ , which although not exactly the same, also  
 331 suggests an exponential growth and slow down rate, and thus coincides with the theory prediction in  
 332 the one-layer model discussed in Appendix A.1. Thus, we prefer to view the theory in this subsection  
 333 as a refinement of the theory derived by one layer model, instead of a refutation. If we take into  
 334 account more layer and non-linearities, we might be able to find more refinements but the predictions  
 335 should follow the same trend as described in this subsection, since all the predicted behaviours of the  
 336 theory presented in this subsection are experimentally verified with more complex models.

## 337 B Model Compositionally Generalize in Topologically Constrained Order

338 In this section, we introduce another phenomenon observed on SIM task learning that we don't put  
 339 in the main paper: the order of compositional generalization happens in a topologically constrained  
 340 order.

341 In this section, instead of the single test point  $\hat{\mathbf{x}}$ , we introduce a hierarchy of test points. Specifically,  
 342 let  $\mathcal{I} = \{0, 1\}^s$  be the index set of test points. For each  $v \in \mathcal{I}$ , we define a test point

$$\hat{\mathbf{x}}^{(v)} = \sum_{p=1}^s v_p \mu_p \mathbf{1}_p, \quad (\text{B.1})$$

<sup>3</sup>And in opposite, when  $s$  is large it's unlikely that all entries are initialized positive, thus the Transient Memorization happens rarer.

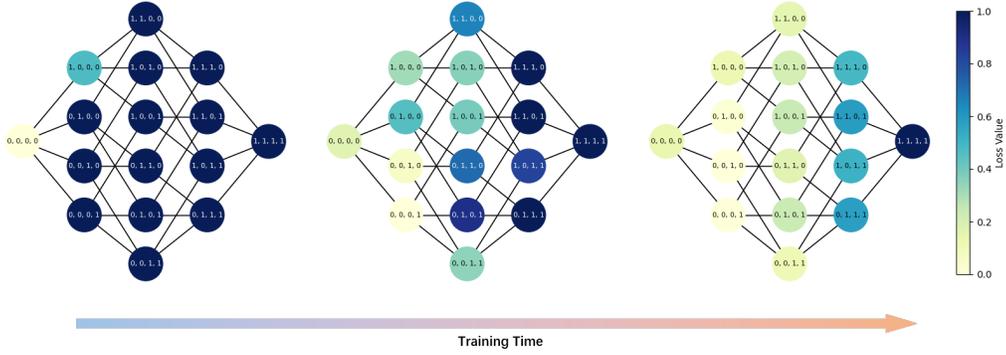


Figure 8: The loss at each test point in different timepoints during training for a 2-layer MLP with ReLU activation. Each graph represents a timepoint. Each node in the graph represents a test point, with index printed on it, and edges connecting test nodes with Hamming distance 1. The color of the graph represents the loss of corresponding test point. Notice that we truncate the loss at 1 in order to unify the scale. From left to right: epoch = 1, 3, 5.

343 and call  $\hat{\mathbf{x}}^{(v)}$  the test point with the index  $v$ . Intuitively, the index  $v$  describes which training sets are  
 344 combined into the current test point. If  $\|v\| = 1$  then  $\hat{\mathbf{x}}^{(v)}$  is the center of one of the training clusters.

345 We assign the component-wise ordering  $\preceq$  to the index set  $\mathcal{I}$ , i.e. for  $\mathbf{u}, \mathbf{v} \in \mathcal{I}$ , we say  $\mathbf{u} \preceq \mathbf{v}$  if and  
 346 only if  $\forall i \in [n], u_i \leq v_i$ . It's easy to see that  $\preceq$  is a partial-ordering.

347 Interestingly, in the SIM experiment, the order of the generalization in different test points strictly  
 348 follow the component-wise order. This finding can be described formally in the following way: the  
 349 loss function is an order homomorphism between  $\preceq$  on the index set, and  $\leq$  on the real number. Let  
 350  $\ell(\mathbf{z})$  be the loss function of the test point  $\mathbf{z}$ , then we have the following empirical observation:

$$\forall \mathbf{u}, \mathbf{v} \in \mathcal{I}, \mathbf{u} \preceq \mathbf{v} \implies \ell(\hat{\mathbf{x}}^{(\mathbf{u})}) \leq \ell(\hat{\mathbf{x}}^{(\mathbf{v})}). \quad (\text{B.2})$$

351 In Figure 8 we show the loss of each test point in several timepoints, with  $\boldsymbol{\mu} = (1, 2, 3, 4)$ ,  $\boldsymbol{\sigma} = \{\frac{1}{2}\}^4$ .  
 352 There is a clear trend that the test points that are on the right of the graph (larger in the component-  
 353 wise order) will only be learned after all of its predecessors are all learned. We call this phenomenon  
 354 the *topological constraint* since the constraint is based on the topology of the graph in Figure 8.

## 355 C Proofs and Calculations

356 In the main text we have omitted some critical proofs and calculations due to space limitation. In this  
 357 section we provide the complete derivations. Notice that we delay the calculations of Appendix A.2  
 358 to Appendix D due to its length.

359 **C.1 The Loss Function with Linear Model and Infinite Data Limit**

360 In this subsection we derive the transformed loss function eq. (A.1), as well as the expression of the  
 361 data matrix  $\mathbf{A}$ . For convenience we denote  $\mathbf{W}_\theta$  by  $\mathbf{W}$ . We have

$$\mathcal{L}(\theta) = \frac{1}{2ns} \sum_{p=1}^s \sum_{k=1}^n \left\| (\mathbf{W} - \mathbf{I}) \mathbf{x}_k^{(p)} \right\|^2 \quad (\text{C.1})$$

$$= \frac{1}{2ns} \text{Tr} \left[ \mathbf{x}_k^{(p)\top} (\mathbf{W} - \mathbf{I})^\top (\mathbf{W} - \mathbf{I}) \mathbf{x}_k^{(p)} \right] \quad (\text{C.2})$$

$$= \frac{1}{2ns} \text{Tr} \left[ (\mathbf{W} - \mathbf{I})^\top (\mathbf{W} - \mathbf{I}) \mathbf{x}_k^{(p)} \mathbf{x}_k^{(p)\top} \right] \quad (\text{C.3})$$

$$= \frac{1}{2} \text{Tr} \left[ (\mathbf{W} - \mathbf{I})^\top (\mathbf{W} - \mathbf{I}) \frac{1}{ns} \mathbf{x}_k^{(p)} \mathbf{x}_k^{(p)\top} \right] \quad (\text{C.4})$$

$$= \frac{1}{2} \text{Tr} \left[ \mathbf{A}^{1/2} (\mathbf{W} - \mathbf{I})^\top (\mathbf{W} - \mathbf{I}) \mathbf{A}^{1/2} \right] \quad (\text{C.5})$$

$$= \frac{1}{2} \left\| (\mathbf{W} - \mathbf{I}) \mathbf{A}^{1/2} \right\|_{\mathcal{F}}^2. \quad (\text{C.6})$$

362 Let  $\mathcal{G}$  be the data generating process. It can be viewed as two components: first assign one of the  $s$   
 363 clusters, and then draw a Gaussian vector from a Gaussian distribution in that cluster. Specifically, let  
 364  $\mathbf{x}$  be an arbitrary sample from the training set, then the distribution of  $\mathbf{x}$  is equal to

$$\mathbf{x} \simeq \boldsymbol{\mu}^{(\eta)} + \text{diag}(\boldsymbol{\sigma}) \boldsymbol{\xi}, \quad (\text{C.7})$$

365 where  $\eta$  is a uniform random variable taking values in  $[s]$  and  $\boldsymbol{\xi} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  is a random Gaussian  
 366 vector that is independent from  $\eta$ . Here  $\simeq$  represents having the same distribution.

367 When  $n \rightarrow \infty$ , the data matrix  $\mathbf{A}$  converges to the true covariance, which is is

$$\mathbf{A} \rightarrow \mathbb{E}(\mathbf{x} \mathbf{x}^\top) \quad (\text{C.8})$$

$$= \mathbb{E} \left[ \left( \boldsymbol{\mu}^{(\eta)} + \text{diag}(\boldsymbol{\sigma}) \boldsymbol{\xi} \right) \left( \boldsymbol{\mu}^{(\eta)} + \text{diag}(\boldsymbol{\sigma}) \boldsymbol{\xi} \right)^\top \right] \quad (\text{C.9})$$

$$= \mathbb{E} \left( \boldsymbol{\mu}^{(\eta)} \boldsymbol{\mu}^{(\eta)\top} \right) + \mathbb{E} \text{diag}(\boldsymbol{\sigma}) \boldsymbol{\xi} \boldsymbol{\xi}^\top \text{diag}(\boldsymbol{\sigma}) \quad (\text{C.10})$$

$$= \frac{1}{s} \sum_{p=1}^s \boldsymbol{\mu}^{(p)} \boldsymbol{\mu}^{(p)\top} + \text{diag}(\boldsymbol{\sigma})^2 \quad (\text{C.11})$$

$$= \frac{1}{s} \sum_{p=1}^s \mu_p^2 \mathbf{1}_p \mathbf{1}_p^\top + \text{diag}(\boldsymbol{\sigma})^2 \quad (\text{C.12})$$

$$= \frac{1}{s} \text{diag}(\boldsymbol{\mu})^2 + \text{diag}(\boldsymbol{\sigma})^2. \quad (\text{C.13})$$

368 **C.2 Proof of Theorem A.1**

369 In this subsection for the notation-wise convenience we denote  $\mathbf{W} = \boldsymbol{\theta}$ . Since the model is one layer,  
 370 the loss function eq. (A.1) becomes

$$\mathcal{L}(\mathbf{W}) = \frac{1}{2} \left\| (\mathbf{W} - \mathbf{I}) \mathbf{A}^{1/2} \right\|_{\mathcal{F}}^2, \quad (\text{C.14})$$

371 and the gradient is

$$\nabla \mathcal{L}(\mathbf{W}) = (\mathbf{W} - \mathbf{I}) \mathbf{A} = \mathbf{W} \mathbf{A} - \mathbf{A}. \quad (\text{C.15})$$

372 We denote the  $k$ -th row of  $\mathbf{W}$  and  $\mathbf{A}$  by  $\mathbf{w}_k$  and  $\mathbf{A}_k$  respectively. Then we have

$$\dot{\mathbf{w}}_k = -\mathbf{A} \mathbf{w}_k + \mathbf{a}_k. \quad (\text{C.16})$$

373 The solution of this differential equation is

$$\mathbf{w}_k(t) = \exp(-\mathbf{A}t) \left[ \mathbf{w}_k(0) - \mathbf{A}^{-1} \mathbf{a}_k \right] + \mathbf{A}^{-1} \mathbf{a}_k, \quad (\text{C.17})$$

374 where we use the convention  $0 \times (0^{-1}) = 0$  to avoid the non-invertible case of  $\mathbf{A}$ .

375 Thus for any  $\mathbf{z} \in \mathbb{R}^d$  we have

$$f(\mathbf{W}(t); \mathbf{z})_k = \langle \mathbf{w}_k(t), \mathbf{z} \rangle \quad (\text{C.18})$$

$$= \langle (\mathbf{I} - e^{-\mathbf{A}t}) \mathbf{A}^{-1} \mathbf{a}_k, \mathbf{z} \rangle + \langle e^{-\mathbf{A}t} \mathbf{w}_k(0), \mathbf{z} \rangle \quad (\text{C.19})$$

$$= \sum_{p=1}^n \frac{1 - e^{-a_p t}}{a_p} \mathbb{1}_{\{k=p\}} a_p z_p + \sum_{i=1}^n e^{-a_i t} w_{k,i}(0) z_i \quad (\text{C.20})$$

$$= \mathbb{1}_{\{k \leq s\}} (1 - e^{-a_k t}) z_k + \sum_{i=1}^n e^{-a_i t} w_{k,i}(0) z_i, \quad (\text{C.21})$$

376 and this proves the claim.

## 377 D Theoretical Analysis of the Two Layer Model

378 In this section we provide a detailed analysis of the symmetric two layer model described in Ap-  
379 pendix A.2.

380 In the theory part we frequently consider functions of a variable  $t$  which is explained as time. If  
381 a function  $g(t)$  is a function of time  $t$ , we denote the derivative of  $g$  w.r.t.  $t$  by  $\dot{g}(t) = \left. \frac{dg}{dt} \right|_{t=t}$ .  
382 Moreover, we sometimes omit the argument  $t$ , i.e.  $g$  means  $g(t)$  for any time  $t$ .

383 For a statement  $\phi$ , we define  $\mathbb{1}_{\{\phi\}} = \begin{cases} 1 & \phi \text{ is true} \\ 0 & \phi \text{ is false} \end{cases}$  that indicates the Boolean value of  $\phi$ .

384 In this section we assume a finite step size, i.e.  $\mathbf{W} : \mathbb{N} \rightarrow \mathbb{R}^{d \times d}$  is initialized by  $\mathbf{W}(0)$  and updated  
385 by

$$\frac{\mathbf{W}(t+1) - \mathbf{W}(t)}{\eta} = -\mathbf{U}(t) \nabla \mathcal{L}(\mathbf{U}(t))^\top - \nabla \mathcal{L}(\mathbf{U}(t)) \mathbf{U}(t)^\top \quad (\text{D.1})$$

$$= \mathbf{W}(t) \mathbf{A} + \mathbf{A} \mathbf{W}(t) - \frac{1}{2} [\mathbf{A} \mathbf{W}(t)^2 + \mathbf{W}(t)^2 \mathbf{A} + 2 \mathbf{W}(t) \mathbf{A} \mathbf{W}(t)]. \quad (\text{D.2})$$

386 The update of each entry  $w_{i,j}(t)$  can be decomposed into three terms, as we described in the main  
387 text:

$$\frac{w_{i,j}(t+1) - w_{i,j}(t)}{\eta} = w_{i,j}(t)(a_i + a_j) - \frac{1}{2} \sum_{k=1}^d w_{k,i} w_{k,j} (a_i + a_j + 2a_k) \quad (\text{D.3})$$

$$= \underbrace{w_{i,j}(t)(a_i + a_j)}_{G_{i,j}(t)} \quad (\text{D.4})$$

$$- \frac{1}{2} w_{i,j} \underbrace{[w_{i,i}(3a_i + a_j) + \mathbb{1}_{\{i \neq j\}} w_{j,j}(3a_j + a_i)]}_{S_{i,j}(t)} \quad (\text{D.5})$$

$$- \frac{1}{2} \underbrace{\sum_{\substack{k \neq i \\ k \neq j}} w_{k,i}(t) w_{k,j}(t) (a_i + a_j + 2a_k)}_{N_{i,j}(t)}. \quad (\text{D.6})$$

### 388 D.1 Assumptions

389 We need make several assumptions to prove the results. Below we make several assumptions that  
390 all commonly hold in the practice. The first assumption to make is that both the value of  $a_k$  and the  
391 initialization of  $\mathbf{W}$  is bounded.

392 **Assumption D.1** (Bounded Initialization and Signal Strength). *There exists  $\alpha > 0, \gamma > 1, \beta > 1$*   
 393 *such that*

$$\begin{aligned} \forall k, \alpha \leq a_k \leq \gamma\alpha, & \quad (D.7) \\ \forall i, j, \omega \leq |w_{i,j}(0)| \leq \beta\omega. & \quad (D.8) \end{aligned}$$

394 The second assumption is that the step size is small enough.

395 **Assumption D.2** (Small Step Size). *There exists a constant  $K \geq 20$ , such that  $\eta \leq \frac{1}{9K\gamma\alpha}$ .*

396 Next, we define a concept called initial phase. The definition of initial phase is related to a constant  
 397  $P > 0$ .

398 **Definition D.1.** *Assume there is a constant  $P > 0$ . For an entry  $(i, j)$  and time  $t$ , if  $|w_{i,j}(t)| \leq P\beta\omega$ ,*  
 399 *we say this entry is in **initial phase**.*

400 The next assumption to make the that the boundary of the initial phase should not be too large.

401 **Assumption D.3** (Small Initial Phase).  $P\omega\beta \leq 0.4$ .

402 The next assumption to make are that the intialization value ( $\omega$ ) should not be too large.

**Assumption D.4** (Small Initialization).

$$\omega \leq \min \left\{ \frac{\min\{\kappa - 1, 1 - \kappa^{-1/2}\}}{PK\gamma d\beta^2}, \frac{1}{\sqrt{2\beta}} \right\} \quad (D.9)$$

403 *and  $\kappa > 1.1$ , and  $\kappa \leq 1 + \frac{1}{2}KC^{-1}$ ,  $P \geq 2$ .*

404 Finally, we also assume that the signal strength difference is significant enough.

405 **Assumption D.5** (Significant Signal Strength Difference). *For any  $i > j$ , we have*

$$\frac{a_i + a_j}{2a_i} \leq \frac{\log P}{10\kappa^2 \log \frac{1}{P\beta\omega} + \log P\beta}. \quad (D.10)$$

406 *and there exists a constant  $C > 1$  such that  $a_i - 3a_j \geq C^{-1}\alpha$ .*

## 407 D.2 The Characterization of the Evolution of the Jacobian

408 In this subsection, we provide a series of lemmas that characterize each stage the evolution of the  
 409 Jacobian matrix  $\mathbf{W}$ .

410 The whole proof is based on induction, and in order to avoid a too complicated induction, we make  
 411 the following assertion, which obviously holds at initialization.

412 **Assertion D.1.** *For all  $t \in \mathbb{N}$ , if  $i \neq j$ , then the entry  $(i, j)$  stays in the initial phase for all time.*

413 We will use Assertion D.1 as an assumption throughout the proves and prove it at the end. This is  
 414 essentially another way of writing inductions.

415 We have the following corollary that directly followed by Assertion D.1.

416 **Corollary D.1.** *For all  $t \in \mathbb{N}$  and all  $i, j$ ,  $|N_{i,j}(t)| \leq 2P\gamma\alpha d\beta^2\omega^2$ .*

417 Now, we are ready to present and prove the major lemmas. The first lemma is to post a (rather loose)  
 418 upper bound of the value of the entries.

419 **Lemma D.1** (Upper Bounded Growth). *Consider entry  $(i, j)$ . We have for all  $t \in \mathbb{N}$ , at timepoint  $t$*   
 420 *the absolute value of the  $(i, j)$ -th entry satisfies*

$$|w_{i,j}(t)| \leq |w_{i,j}(0)| \exp[\eta t(a_i + a_j)\kappa]. \quad (D.11)$$

421 *Proof.* Since of the  $N_{i,j}$  term we only use its absolute value, the positive case and negative case are  
 422 symmetric. WLOG we only consider the case where  $w_{i,j}(0) > 0$  here.

423 The claim is obviously satisfied at initialization. We use it as the inductive hypothesis. Suppose at  
 424 timepoint  $t \leq T - 1$  the claim is satisfied, we consider the time step  $t + 1$ .

425 Since Assertion D.1 guaranteed that every non-diagonal entry is in the initial phase, and the  $S_{i,j}$  term  
 426 has different symbol with  $w_{i,j}(0)$ , we have

$$S_{i,j}(t) + N_{i,j}(t) \leq 2P\gamma\alpha d\beta^2\omega^2. \quad (\text{D.12})$$

427 We have

$$w_{i,j}(t+1) - w_{i,j}(t) \leq \eta w_{i,j}(t)(a_i + a_j) + 4\eta\gamma\alpha d\beta_0\omega^2 \quad (\text{D.13})$$

$$\leq \eta(a_i + a_j)w_{i,j}(0) \exp[\eta t(a_i + a_j)\kappa] + 2P\eta\gamma\alpha d\beta^2\omega^2 \quad (\text{D.14})$$

$$= w_{i,j}(0) \exp[\eta t(a_i + a_j)\kappa] \left[ \eta(a_i + a_j) + \frac{2P\eta\gamma\alpha d\beta^2\omega^2}{w_{i,j}(0) \exp[\eta t(a_i + a_j)\kappa]} \right] \quad (\text{D.15})$$

428 From Assumption D.4, we have

$$\eta(a_i + a_j) + \frac{2P\eta\gamma\alpha d\beta^2\omega^2}{w_{i,j}(0) \exp[\eta t(a_i + a_j)\kappa]} \leq \eta(a_i + a_j) + 2P\eta\gamma\alpha d\beta^2\omega \quad (\text{D.16})$$

$$\leq \eta(a_i + a_j) + 2(\kappa - 1)\eta\alpha \quad (\text{D.17})$$

$$\leq \kappa\eta(a_i + a_j) \quad (\text{D.18})$$

$$\leq \exp(\kappa\eta[a_i + a_j]) - 1, \quad (\text{D.19})$$

429 thus we have

$$w_{i,j}(t+1) \leq w_{i,j}(t) + [\exp(\kappa\eta[a_i + a_j]) - 1] w_{i,j}(t) \quad (\text{D.20})$$

$$\leq w_{i,j}(0) \exp[\eta(t+1)(a_i + a_j)\kappa]. \quad (\text{D.21})$$

430 Finally, notice that since  $T_1 = \frac{\kappa \log P}{2\eta\gamma\alpha} \leq \frac{\kappa \log 2}{\eta(a_i + a_j)}$ , we have

$$\exp[\eta T_1(a_i + a_j)\kappa^{-1}] \leq P. \quad (\text{D.22})$$

431

□

432 Next, we prove that Lemma D.1 is tight in the initial stage of the training, up to a constant  $\kappa$  in the  
 433 exponential term.

434 **Lemma D.2** (Lower Bounded Initial Growth). *Let  $T_1 = \frac{\log P}{2\eta\gamma\alpha\kappa}$ . We have for all  $t \in [T_1]$ , at timepoint  
 435  $t$  every entry  $(i, j)$  is in the initial phase, and the absolute value of the  $(i, j)$ -th entry satisfies*

$$|w_{i,j}(t)| \geq |w_{i,j}(0)| \exp[\eta t(a_i + a_j)\kappa^{-1}] \quad (\text{D.23})$$

436 and  $w_{i,j}(t)w_{i,j}(0) > 0$ .

437 *Proof.* Similar to the proof of Lemma D.1, we may just assume  $w_{i,j}(0) > 0$ .

438 Moreover, we also use the claim as an inductive hypothesis and prove it by induction. Since here the  
 439 inductive hypothesis states that every entry is in the initial phase, we have

$$|S_{i,j}(t) + N_{i,j}(t)| \leq 4\gamma\alpha d\beta^2\omega^2. \quad (\text{D.24})$$

440 We have

$$w_{i,j}(t+1) - w_{i,j}(t) \geq \eta(a_i + a_j)w_{i,j}(0) \exp[\eta t(a_i + a_j)\kappa^{-1}] - 2P\eta\gamma\alpha d\beta^2\omega^2 \quad (\text{D.25})$$

$$= w_{i,j}(0) \exp[\eta t(a_i + a_j)\kappa^{-1}] \left[ \eta(a_i + a_j) - \frac{2P\eta\gamma\alpha d\beta^2\omega^2}{w_{i,j}(0) \exp[\eta t(a_i + a_j)\kappa^{-1}]} \right] \quad (\text{D.26})$$

441 From Assumption D.4, we have

$$\frac{2P\eta\gamma\alpha d\beta^2\omega^2}{w_{i,j}(0) \exp[\eta t(a_i + a_j)\kappa^{-1}]} \leq 2P\eta\gamma\alpha d\beta^2\omega \quad (\text{D.27})$$

$$\leq \left(1 - \kappa^{-1/2}\right) \eta(a_i + a_j). \quad (\text{D.28})$$

442 Moreover, notice that when  $\kappa > 1.1$ , for any  $x < 0.1$ , we have  $\kappa^{-1/2}x + 1 \geq e^{\kappa^{-1}x}$ . Since  
 443 Assumption D.2 ensured that  $\eta \leq \frac{1}{10(a_i + a_j)}$ , we have

$$w_{i,j}(t+1) \geq w_{i,j}(t) + w_{i,j}(t) \left[ \kappa^{-1/2} \eta (a_i + a_j) \right] \quad (\text{D.29})$$

$$\geq w_{i,j}(t) \exp(\eta(a_i + a_j)\kappa^{-1}) \quad (\text{D.30})$$

$$\geq w_{i,j}(0) \exp[\eta(t+1)(a_i + a_j)\kappa^{-1}]. \quad (\text{D.31})$$

444 Finally, from lemma D.1, we have when

$$w_{i,j}(t) \leq |w_{i,j}(0)| \exp(\eta t(a_i + a_j)\kappa) \quad (\text{D.32})$$

$$\leq \beta\omega \exp(2\eta T_1 \gamma \alpha \kappa) \quad (\text{D.33})$$

$$\leq P\beta\omega, \quad (\text{D.34})$$

445 which confirms that every entry  $(i, j)$  stays in the initial phase before time  $T_1$ .

446

□

447 Notice that the time bound in Lemma D.2 is a uniform one which applies to all entries. For the major  
 448 entries, we might want to consider a finer bound of the time that it leaves the initial phase. This can  
 449 be proved by essentially repeating the same proof idea of Lemma D.2.

450 **Lemma D.3** (Lower Bounded Initial Growth for Diagonal Entries). *Consider an diagonal entry  $(i, i)$ .*

451 *Let  $T_1^{(i)} = \frac{\log \frac{P\beta\omega}{w_{i,i}(0)}}{2\eta a_i \kappa}$ . We have for all  $t \in [T_1^{(i)}]$ , at timepoint  $t$  the entry  $(i, i)$  is in the initial phase,  
 452 and the absolute value of the  $(i, i)$ -th entry satisfies*

$$w_{i,i}(t) \geq w_{i,i}(0) \exp[2\eta t a_i \kappa^{-1}]. \quad (\text{D.35})$$

453 We omit the proof of Lemma D.3 since it is almost identical to the proof of Lemma D.2, only with  
 454 replacing  $\gamma\alpha$  by  $a_i$  and  $\beta\omega$  by  $w_{i,i}(0)$ .

455 Next, we characterize the behavior of one diagonal entry after it leaves the initial phase.

456 **Lemma D.4** (Lower Bounded After-Initial Growth for Diagonal Entries). *Consider a diagonal entry  
 457  $(i, i)$ . If at time  $t_0$  we have  $|w_{i,i}(t_0)| \geq P\beta\omega$ , and for a  $\lambda \in (P\beta\omega, 1 - K^{-1})$ , before time  $T^{(\lambda)}$  we  
 458 have  $w_{i,i}(t + t_0) < \lambda$  for all  $t \in [T^{(\lambda)}]$ , then we have*

$$w_{i,i}(t + t_0) \geq w_{i,i}(t_0) \exp[2\eta t a_i (1 - \lambda)\kappa^{-1}]. \quad (\text{D.36})$$

459 *Moreover,  $w_{i,i}(0), w_{i,i}(t_0), w_{i,i}(t_0 + t) \geq 0$ .*

460 *Proof.* Notice that since  $\mathbf{W} = \mathbf{U}\mathbf{U}^\top$  is a PSD matrix, its diagonal entries are always non-negative,  
 461 this ensures that  $w_{i,i}(0), w_{i,i}(t_0), w_{i,i}(t_0 + t) \geq 0$ .

462 For the time after  $t_0$  and before  $t_0 + T^{(\lambda)}$ , we use an induction to prove the claim, with the claim  
 463 itself as the inductive hypothesis. It clearly holds when  $t = 1$ .

464 Notice that when  $w_{i,j}(t') < \lambda$ , we have

$$G_{i,j}(t') - S_{i,j}(t') = 2a_i w_{i,i}(t') [1 - w_{i,i}(t')] \geq 2a_i w_{i,i}(t') (1 - \lambda). \quad (\text{D.37})$$

465 Thus we have

$$w_{i,i}(t_0 + t + 1) - w_{i,i}(t_0 + t) \quad (\text{D.38})$$

$$\geq 2\eta a_i (1 - \lambda) w_{i,i}(t_0) \exp[\eta t (a_i + a_j)(1 - \lambda)\kappa^{-1}] - 2P\eta\gamma\alpha d\beta^2\omega^2 \quad (\text{D.39})$$

$$= w_{i,i}(t_0) \exp[2\eta t a_i (1 - \lambda)\kappa^{-1}] \left[ 2\eta a_i (1 - \lambda) - \frac{2P\eta\gamma\alpha d\beta^2\omega^2}{w_{i,i}(t_0) \exp[2\eta t a_i (1 - \lambda)\kappa^{-1}]} \right] \quad (\text{D.40})$$

466 Since  $\lambda < 1 - K^{-1}$ , and  $w_{i,i}(t_0) \geq 2\beta\omega \geq \omega$ , from Assumption D.4, we have

$$\frac{2P\eta\gamma\alpha d\beta^2\omega^2}{w_{i,i}(t_0) \exp[2\eta t a_i (1 - \lambda)\kappa^{-1}]} \leq 2P\eta\gamma\alpha d\beta^2\omega \quad (\text{D.41})$$

$$\leq 2K^{-1} \left(1 - \kappa^{-1/2}\right) \eta\alpha \quad (\text{D.42})$$

$$\leq 2 \left(1 - \kappa^{-1/2}\right) \eta a_i (1 - \lambda). \quad (\text{D.43})$$

467 Moreover, since Assumption D.2 ensured that  $\eta \leq \frac{1}{2Ka_i(1-\lambda)} \leq \frac{1}{20a_i(1-\lambda)}$ , using the fact that if  
 468  $\kappa > 1.1$  then  $\kappa^{-1/2}x + 1 \geq e^{\kappa^{-1}x}$  for any  $x < 0.1$ , we can get

$$w_{i,i}(t+1) \geq w_{i,i}(t) + w_{i,i}(t) \left[ \kappa^{-1/2} 2\eta a_i (1-\lambda) \right] \quad (\text{D.44})$$

$$\geq w_{i,i}(t) \exp \left( 2\eta a_i \kappa^{-1} (1-\lambda) \right) \quad (\text{D.45})$$

$$\geq w_{i,i}(t_0) \exp \left[ 2\eta (t+1) \kappa^{-1} (1-\lambda) \right]. \quad (\text{D.46})$$

469

□

470 Next, we provide an uniform upper bound (over time) of the diagonal entries. Remember that we  
 471 mentioned in the gradient flow case, the diagonal term stops evolving when it reaches 1. In the  
 472 discrete case, since the step size is not infinitesimal, Lemma D.5 shows that it can actually exceed 1 a  
 473 little bit but not too much since the step size is small.

474 **Lemma D.5** (Upper Bounded Diagonal Entry). *For any diagonal entry  $(i, i)$  and any time  $t$ ,  $0 \leq$   
 475  $w_{i,i}(t) \leq 1 + 2K^{-1}$ .*

476 *Proof.* First notice that since  $\mathbf{W}(t)$  is PSD, its diagonal entry  $w_{i,i}(t)$  should always be non-negative,  
 477 thus  $w_{i,i}(t) \geq 0$  is always satisfied. In the following we prove  $w_{i,i}(t) \leq 1 + 2K^{-1}$ .

478 We use induction to prove this claim. The inductive hypothesis is the claim it self. It is obviously  
 479 satisfied at initialization. In the following we assume the claim is satisfied at timepoint  $t$  and prove it  
 480 for timepoint  $t+1$ . Notice that since  $K \leq 10$ , we have  $1 + K^{-1} \leq 2$ .

481 Notice that by Assertion D.1 and Assumption D.4,

$$|N_{i,i}(t)| \leq 2P\gamma\alpha d\beta^2\omega^2 \leq \frac{(\kappa-1)^2}{K^2\gamma d\beta^2}\alpha \leq K^{-1}a_i. \quad (\text{D.47})$$

482 If  $w_{i,i}(t) \geq 1 + K^{-1}$ , we have

$$G_{i,i}(t) - S_{i,i}(t) = 2a_i w_{i,i}(t)(1 - w_{i,i}(t)) \leq -4a_i K^{-1}. \quad (\text{D.48})$$

483 Therefore,

$$w_{i,i}(t+1) = w_{i,i}(t) + \eta [G_{i,i}(t) - S_{i,i}(t) - N_{i,i}(t)] \quad (\text{D.49})$$

$$\leq w_{i,i}(t) - 3a_i K^{-1} \eta \quad (\text{D.50})$$

$$\leq w_{i,i}(t) \quad (\text{D.51})$$

$$\leq 1 + 2K^{-1}. \quad (\text{D.52})$$

484 Moreover, since  $w_{i,i}(t) \leq 1 + 2K^{-1} \leq 2$ , we have

$$|G_{i,i}(t)| + |S_{i,i}(t)| + |N_{i,i}(t)| \leq 4a_i + 4a_i + K^{-1}a_i \leq 9\gamma\alpha \leq \frac{1}{K\eta}. \quad (\text{D.53})$$

485 When  $w_{i,i}(t) \leq 1 + K^{-1}$ , we have

$$w_{i,i}(t+1) \leq w_{i,i}(t) + \eta (|G_{i,i}(t)| + |S_{i,i}(t)| + |N_{i,i}(t)|) \leq 1 + 2K^{-1}. \quad (\text{D.54})$$

486 The above results together shows that  $w_{i,i}(t+1) \leq 1 + 2K^{-1}$ .

487

□

488 **Corollary D.2** (Upper Bounded Diagonal Update). *For any diagonal entry  $(i, i)$  and any time  $t$ ,*  
 489  $|w_{i,i}(t+1) - w_{i,i}(t)| \leq K^{-1}$ .

490 Corollary D.2 is a direct consequence of Lemma D.5 (and we actually proved Corollary D.2 in the  
 491 proof of Lemma D.5).

492 The next lemma lower bounds the final value of diagonal entries. Together with Lemma D.5 we  
 493 show that in the terminal stage of training the diagonal entries oscillate around 1 by the amplitude not  
 494 exceeding  $2K^{-1}$ .

495 **Lemma D.6.** Consider a diagonal entry  $(i, i)$ . If at time  $t_0$  we have  $w_{i,i}(t_0) \geq 1 - 2K^{-1}$ , then for  
 496 all  $t' \geq t_0$  we have  $w_{i,i}(t') \geq 1 - 2K^{-1}$ .

497 *Proof.* we use an induction. The inductive hypothesis the claim itself. This obviously holds when  
 498  $t' = t_0$ . We assume  $w_{i,i}(t') \geq 1 - 2K^{-1}$  at timepoint  $t'$  and prove the claim for  $t' + 1$ .

499 If  $w_{i,i}(t') < 1 - K^{-1}$ , then from Lemma D.4 we know

$$w_{i,i}(t' + 1) \geq w_{i,i}(t') \geq 1 - 2K^{-1}. \quad (\text{D.55})$$

500 If  $w_{i,i}(t') > 1 - K^{-1}$ , then from Corollary D.2 we have

$$w_{i,i}(t' + 1) \geq w_{i,i}(t') - K^{-1} \geq 1 - 2K^{-1}. \quad (\text{D.56})$$

501  $\square$

502 Now, we are ready to prove Assertion D.1 by considering the suppression. We first prove a lemma  
 503 that upper bounds the absolute value of the minor entries after its corresponding major entry becomes  
 504 significant.

505 **Lemma D.7 (Suppression).** Consider an off-diagonal entry  $(i, j)$  where  $i > j$ . If there exists a time  
 506  $t_0$  such that  $w_{i,i}(t_0) > 0.8$ , then for any  $t' \geq t_0$  we have

$$|w_{i,j}(t')| \leq \max \{|w_{i,j}(t_0)|, \omega\}. \quad (\text{D.57})$$

507 *Proof.* Since  $K > 10$ , from Lemma D.6 and Lemma D.4 we know  $w_{i,i}(t') > 0.8$  for all  $t' \geq t_0$ .

508 In this proof, we use an induction with the inductive hypothesis being the claim itself, i.e. we assume  
 509 the claim is true at timepoint  $t'$  and prove it for  $t' + 1$ . The claim obviously holds for  $t' = t_0$ .

510 Since in this proof we only use the absolute value of  $N_{i,j}$ , WLOG we may assume that  $w_{i,j}(t') > 0$ .

511 If  $w_{i,j}(t') < \omega$  then we have proved the claim. In the following we may assume  $w_{i,j}(t') \geq \omega$ .

512 We have

$$G_{i,j}(t') - S_{i,j}(t') \leq w_{i,j}(t')(a_i + a_j) - \frac{1}{2}w_{i,j}(t')w_{i,i}(3a_i + a_j) \quad (\text{D.58})$$

$$\leq w_{i,j}(t')(a_i + a_j) - w_{i,j}(t')[0.4(3a_i + a_j)] \quad (\text{D.59})$$

$$= -\frac{1}{5}w_{i,j}(t')a_i + \frac{3}{5}w_{i,j}(t')a_j \quad (\text{D.60})$$

$$\stackrel{(i)}{\leq} -C^{-1}\omega\alpha, \quad (\text{D.61})$$

513 where in (i) we use Assumption D.5.

514 Thus we have

$$G_{i,j}(t') - S_{i,j}(t') - N_{i,j}(t') \leq G_{i,j}(t') - S_{i,j}(t') + |N_{i,j}(t')| \quad (\text{D.62})$$

$$\leq -C^{-1}\omega\alpha + 2P\gamma\alpha d\beta^2\omega^2 \quad (\text{D.63})$$

$$\stackrel{(i)}{<} 0, \quad (\text{D.64})$$

515 where (i) is from Assumption D.4 and Assumption D.5. This confirms that  $w_{i,j}(t' + 1) < w_{i,j}(t') \leq$   
 516  $\max \{|w_{i,j}(t_0)|, \omega\}$ .

517 Next, we prove  $w_{i,j}(t' + 1) \geq -\max \{|w_{i,j}(t_0)|, \omega\}$ . Notice that Lemma D.5 stated that  $|w_{i,i}| \leq 2$ .

518 Notice that we also have  $w_{i,j}(t') \leq K^{-1}$ , thus

$$|G_{i,j}(t')| + |S_{i,j}(t')| + |N_{i,j}(t')| \leq 10\gamma\alpha|w_{i,j}(t')| + 2P\gamma\alpha d\beta^2\omega^2 \quad (\text{D.65})$$

$$\leq \frac{10|w_{i,j}(t')| + 2Pd\beta^2\omega^2}{9K\eta} \quad (\text{D.66})$$

$$\leq \frac{10|w_{i,j}(t')| + 2\omega}{9K\eta} \quad (\text{D.67})$$

$$\leq \frac{|w_{i,j}(t')| + \omega}{2\eta}. \quad (\text{D.68})$$

519 We have

$$w_{i,j}(t' + 1) \geq w_{i,j}(t') - \eta(|G_{i,j}(t')| + |S_{i,j}(t')| + |N_{i,j}(t')'|) \quad (\text{D.69})$$

$$\geq -\eta(|G_{i,j}(t')| + |S_{i,j}(t')| + |N_{i,j}(t')'|) \quad (\text{D.70})$$

$$\geq -\frac{1}{2}(|w_{i,j}(t')| + \omega) \quad (\text{D.71})$$

$$\geq -\max\{|w_{i,j}(t')|, \omega\}. \quad (\text{D.72})$$

520

□

521 With all the lemmas proved above, we are now ready to prove Assertion D.1.

522 **Lemma D.8** (Assertion D.1). *For all  $t \in \mathbb{N}$ , if  $i \neq j$ , then the entry  $(i, j)$  stays in the initial phase*  
 523 *for all time.*

524 *Proof.* Notice that since  $\mathbf{W}$  is symmetric, we only need to prove the claim for  $i > j$ . Moreover,  
 525 From Lemma D.7, we only need to prove that there exists a timepoint  $t^*$ , such that  $w_{i,i}(t^*) \geq 0.8$ ,  
 526 and  $|w_{i,j}(t^*)| \leq P\beta\omega$ .

527 Let  $t_0 = \frac{\log \frac{P\beta\omega}{w_{i,i}(0)}}{2\eta a_i \kappa}$ , by Lemma D.3, we have  $w_{i,i}(t_0) \geq P\beta\omega$ . By Lemma D.3 and Lemma D.4, we  
 528 have for any  $t \geq t_0$  such that  $w_{i,i}(t) \leq \lambda$ , where  $\lambda = 0.85$ ,

$$w_{i,i}(t) \geq w_{i,i}(t_0) \exp[0.3\eta(t - t_0)a_i\kappa^{-1}] \quad (\text{D.73})$$

$$\geq P\beta\omega \exp[0.3\eta(t - t_0)a_i\kappa^{-1}] \quad (\text{D.74})$$

529 Let  $t'$  be the first time that  $w_{i,i}(t')$  arrives above 0.8. Let  $t^* = \min\left\{\frac{\kappa \log \frac{0.8}{P\beta\omega}}{0.3\eta a_i} + t_0, t'\right\} \geq t_0$ . If

530  $t^* = t'$ , we have  $w_{i,i}(t^*) \geq 0.8$ . If  $t^* = \frac{\log \frac{P\beta\omega}{w_{i,i}(0)}}{2\eta a_i \kappa} + t_0$ , we have

$$w_{i,i}(t^*) \geq w_{i,i}(0) \exp(0.3\eta t^* a_i \kappa^{-1}) \quad (\text{D.75})$$

$$\geq P\beta\omega \exp\left(\log \frac{0.8}{P\beta\omega}\right) \quad (\text{D.76})$$

$$\geq 0.8. \quad (\text{D.77})$$

531 Moreover, from Lemma D.1 and Assumption D.5, we have

$$|w_{i,j}(t^*)| \leq |w_{i,j}(0)| \exp[\eta t^* \kappa(a_i + a_j)] \quad (\text{D.78})$$

$$\leq \beta\omega \exp\left[\left(\frac{\kappa^2 \log \frac{0.8}{P\beta\omega}}{0.15} + \log \frac{P\beta\omega}{w_{i,i}(0)}\right) \times \frac{a_i + a_j}{2a_i}\right] \quad (\text{D.79})$$

$$\leq \beta\omega \exp\left[\left(10\kappa^2 \log \frac{1}{P\beta\omega} + \log P\beta\right) \times \frac{a_i + a_j}{2a_i}\right] \quad (\text{D.80})$$

$$\leq \beta\omega \exp[\log(P)] \quad (\text{D.81})$$

$$\leq P\omega\beta. \quad (\text{D.82})$$

532 The claim is thus proved by combining the above bounds on  $|w_{i,j}(t^*)|$  and  $w_{i,i}(t^*)$  with Lemma D.7.

533

□

534 **E Debunking Challenge Submission**

535 **E.1 What commonly-held position or belief are you challenging?**

536 *Provide a short summary of the body of work challenged by your results. Good summaries should*  
537 *outline the state of the literature and be reasonable, e.g. the people working in this area will agree*  
538 *with your overview. You can cite sources beside published work (e.g., blogs, talks, etc).*

539 People generally believe that double descent with respect to training time (or “epochwise double  
540 descent”) only happens either 1) with large step size or noisy training process; or 2) under over-  
541 parameterized settings.

542 **E.2 How are your results in tension with this commonly-held position?**

543 *Detail how your submission challenges the belief described in (1). You may cite or synthesize results*  
544 *(e.g. figures, derivations, etc) from the main body of your submission and/or the literature.*

545 In Section 3.2, we show a epochwise double descent phenomenon that happens when using large  
546 training set and small step size. Our theoretical analysis in Appendix A.2 shows that even in gradient  
547 flow (infinitesimal stepsize) and infinite data limit, this epochwise double descent still exists (see  
548 Figure 7). This suggests that the epochwise double descent found in our setting is not introduced by  
549 noise in the training or overfitting the data but an inherent property of the task itself.

550 We note that this contradictory comes from the out-of-distribution nature of the SIM task we consider:  
551 when training is noiseless and model is under-parameterized, generally there will not be in-distribution  
552 epochwise double descent, but once we consider an out-of-distribution task, even if it is as simple as  
553 learning identity mapping, there can still be epochwise double descent.

554 **E.3 How do you expect your submission to affect future work?**

555 *Perhaps the new understanding you are proposing calls for new experiments or theory in the area, or*  
556 *maybe it casts doubt on a line of research.*

557 We would like to call for the awareness of this kind of epochwise double descent that is intrinsic to  
558 certain OOD tasks. Since SIM is a very simple toy task but still show epochwise double descent, we  
559 hypothesize that this kind of epochwise double descent might presence across many tasks. It is also  
560 an important direction for future work that to (either empirically or theoretically) fully characterize  
561 the scenarios that has this kind of epochwise double descent. Specifically, since it is not caused by  
562 training or (in-distribution) generalization issues, it is likely caused by the structure of the input data.  
563 Therefore, to determine what kind of input data structure will / will not cause this kind of epochwise  
564 double descent is a very interesting direction to explore.