
A Linear Network Theory of Iterated Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Language provides one of the primary examples of human’s ability to system-
2 atically generalize — reasoning about new situations by combining aspects of
3 previous experiences. Consequently modern machine learning has drawn much in-
4 spiration from linguistics. A recent example is iterated learning, a procedure where
5 generations of networks learn from the output of earlier learners. The result is
6 a refinement of the network’s “language” or output labels for given inputs towards
7 compositional structure. Yet, studies of iterated learning and its application to ma-
8 chine learning have remained empirical. Here we theoretically study the emergence
9 of compositional language, and the ability of simple neural networks to leverage
10 this compositionality to systematically generalize. We build on prior theoretical
11 work on linear networks, which mathematically defines systematic generalization,
12 by extending the analysis of shallow and deep linear network learning dynamics
13 to the iterated learning procedure by deriving exact dynamics to the learning over
14 generations. Our results confirm a long standing conjecture: that multiple genera-
15 tions of iterated learning are required for compositional structure to emerge, which
16 can outperform a single generation network trained with optimal early-stopping.
17 Finally, we show that IL requires depth in the network architecture to be effective
18 and that IL is able to extract modules which systematically generalize.

19 1 Introduction

20 Deep learning techniques have made great strides in tasks like machine translation and language
21 prediction, providing proof of principle that they can succeed in quasi-compositional domains. How-
22 ever, these methods are typically data hungry and the same networks often fail to generalize in even
23 simple settings when training data are scarce (Lake & Baroni, 2018; Lake et al., 2019). *Systematically*
24 *generalize*, leveraging specific learning experiences in diverse new settings (Lazaridou et al., 2018;
25 Lake et al., 2019; Ren et al., 2019), has been proposed as a key feature of intelligent learning agents
26 which can efficiently generalize to novel stimuli in their environment (Hockett & Hockett, 1960;
27 Fodor & Pylyshyn, 1988; Hadley, 1993; Kirby et al., 2015; Lake et al., 2017). Empirically, the degree
28 of systematicity in deep networks is influenced by many factors. One possibility is that the learning
29 dynamics in a deep network could impart an implicit inductive bias toward compositional structure
30 (Hupkes et al., 2020); however, a number of studies have identified situations where depth alone is
31 insufficient for structured generalization (Lake & Baroni, 2018; Niklasson & Sharkey, 1992; Pollack,
32 1990; Phillips & Wiles, 1993; Jarvis et al., 2023). Another significant factor is architectural modularity,
33 which can enable a system to generalize when modules are appropriately configured (Vani et al., 2021;
34 Phillips, 1995; Andreas et al., 2016; Hu et al., 2017, 2018). However, identifying the right modularity
35 through learning remains challenging (Bahdanau et al., 2019; Jarvis et al., 2023). A third possibility
36 builds on iterated learning (IL), a method in which generations of agents train briefly on a language
37 produced by their parent, and then generate a new language for their child (Kirby, 2001; Kalish et al.,
38 2007; Kirby et al., 2015; Vani et al., 2021; Lu et al., 2020a,b). If compositional components are

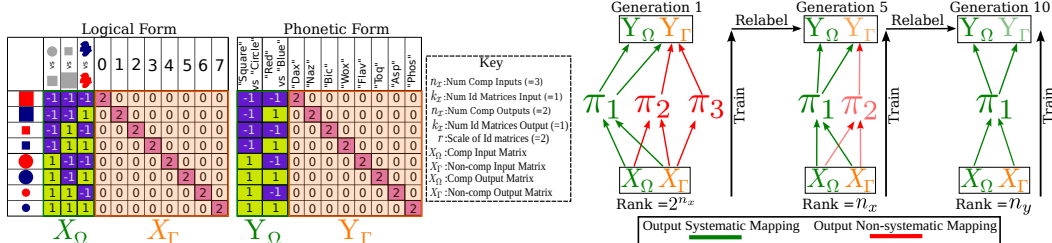


Figure 1: (Left) We schematize the setting with a space of datasets containing compositional (X_Ω) and non-compositional (X_Γ) features in the input (left panel). The task then is to map from these features to a phonetic form for each object which could be by the composition of descriptive words (Y_Ω) for example "the small red square" or by memorizing a name for each objects (Y_Γ) for example "the bic". (Right) The iterated learning procedure: Generations of agents learn from languages generated by their parent, and pass on their acquired language to their children. The portions of the language which are easier to learn are maintained over generations, while difficult language is lost. We demonstrates this process on a linear neural network and prove that IL is able to refine the language to depend on a minimal set of necessary singular values (π_1). This figure also summarizes the notation, structure and singular value decomposition of our space of datasets which is key for the theoretical results which follow.

39 easier to learn than non-compositional ones, this process can successively refine a language towards
 40 compositional structure, which has been hypothesized to be the cause of the compositional nature of
 41 natural language (Kirby, 2001; Kirby et al., 2008). In this work we aim to expand on the deep linear
 42 network framework of Saxe et al. (2019) (Saxe et al., 2019) and the formal analysis of systematicity of
 43 Jarvis et al. (2023) (Jarvis et al., 2023) to analyse the ability of IL to produce compositional language.

44 2 Iterated Learning Dynamics

45 The generalization abilities of deep networks depends on a complex interplay of learning dynamics
 46 (Saxe et al., 2014), architecture (Lake & Baroni, 2018), initialization (Geiger et al., 2020), and dataset
 47 structure (Jarvis et al., 2023). Prior work has demonstrated that gradient descent dynamics alone
 48 do not implicitly favour systematicity (Jarvis et al., 2023) for all but the most compositional of
 49 datasets. To establish whether *generations* of gradient descent learners have an implicit bias towards
 50 systematicity we obtain closed-form learning dynamics for neural networks in the IL procedure. We
 51 build on known exact solutions to the dynamics of learning from small random weights in deep linear
 52 networks (Saxe et al., 2014, 2019) to describe the full learning trajectory analytically.

53 In particular, consider a single hidden layer network computing output $\hat{y} = W^2 W^1 x$ in response
 54 to an input x , trained to minimize the quadratic loss $L(W^1, W^2) = \frac{1}{2} \|Y - W^2 W^1 X\|_2^2$ using full
 55 batch gradient descent. We review the derivation of the linear network dynamics in Appendix A.
 56 However, the main idea is to perform a change of variables so that we track the dynamics of the
 57 network's singular values rather than the individual weights. This does assume the network is feature
 58 learning such that its singular vectors align to those prescribed by the dataset statistics. This is a
 59 reasonable assumption from small initial weights (Saxe et al., 2019) and the singular vectors can
 60 be thought of conceptually as the features learned by the network. The trajectory of each network
 61 effective singular value ($\pi_\alpha(t)$) is described as

$$\pi_\alpha(t) = \frac{\lambda_\alpha / \delta_\alpha}{1 - \left(1 - \frac{\lambda_\alpha}{\delta_\alpha \pi_0}\right) \exp\left(\frac{-2\lambda_\alpha t}{\tau}\right)}. \quad (1)$$

62 These dynamics describe the singular value's trajectory which begins at the initial value π_0 when
 63 $t = 0$ and increases to $\lambda_\alpha / \delta_\alpha$ as $t \rightarrow \infty$. From these dynamics it is helpful to note that the time-
 64 course of the trajectory is only dependent on the input-output covariance matrix (Σ^{y^x}) singular values
 65 (λ_α). Thus, the input covariance (Σ^x) (and its singular values δ_α), affects the stable point of the
 66 network singular values but not the rate of learning.

67 With IL each generation learns from the "language" acquired by the previous generation (Figure 1).
 68 To instantiate this setting, we start from a particular dataset, but halt training before full convergence
 69 after a pre-defined number of training steps. We then use the network's output (logits) as the target
 70 outputs for the next generation. From very early on in training, learning occurs along the modes

71 of variation determined by the dataset statistics. Consequently, the dataset’s singular vectors, and
 72 the features learned by the network, will be maintained for all generations. It is merely the singular
 73 values or salience of the features which changes. Noting this fact permits straightforward analysis
 74 of iterated learning dynamics. Thus, for generation $G > 0$ of learning the asymptote of the network’s
 75 mapping ($\lambda_\alpha^G/\delta_\alpha$) is equal to the effective singular value of the network at the end of the previous
 76 generation of training (π_α^{G-1}). Here λ_α^0 and δ_α^0 are the singular values from the original dataset.
 77 Thus, by a recursive application of Equation 1 we can model the full dynamics of iterated learning:

$$\pi_\alpha^G(t) = \frac{\lambda_\alpha^G/\delta_\alpha}{1 - \left(1 - \frac{\lambda_\alpha^G}{\delta_\alpha \pi_0}\right) \exp\left(\frac{-2\lambda_\alpha^G}{\tau} t\right)} = \frac{\pi_\alpha^{G-1}}{1 - \left(1 - \frac{\pi_\alpha^{G-1}}{\pi_0}\right) \exp\left(\frac{-2\lambda_\alpha^G}{\tau} t\right)} \quad (2)$$

78

79 3 Theoretical Results

79

80 3.1 The Requirement of Multiple Generations

81 By ending a generation of training before the convergence of some effective singular values we will
 82 be decreasing the input-output singular values (λ_α) for the next generation. However, since these
 83 singular values also determine how quickly the mode is learned, this also means it will be learned
 84 slower for subsequent generations. The result is that it is eroded until it is removed from the “language”
 85 (has a singular value of 0). Since our analysis in this work relies on the relative speed at which modes
 86 start and finish learning it is helpful to introduce two new terminology. Firstly, by “escaping time”
 87 (denoted by \hat{t}_α) we refer to the time taken for a mode (indexed by α) to grow meaningfully larger than
 88 0: $\pi_\alpha > \rho$ for a small value of ρ . Secondly, “hitting time” (denoted by t_α^*) refers to the time taken for
 89 a mode (indexed by α) to converge to its final value: $\pi_\alpha - (\lambda_\alpha/\delta_\alpha) < \rho$. We derive explicit equations
 90 for \hat{t}_α and t_α^* in terms of the dataset singular values using Equation 2 in Appendix B. The equations
 91 themselves do not offer immediate insight beyond the concept they represent and so are omitted here.
 92 They are, however necessary for the proves of the theorem and observations which follow.

93 A lingering question in the use of IL has been whether multiple generations of learning is actually
 94 necessary. This is in contrast to a hypothetical optimal early stopping point which would provide all
 95 the same benefits as IL but within a single generation. It is important to note that IL must maintain
 96 all information which we do not wish to remove (maintained modes). There is naturally a trivial
 97 early-stopping time which removes all modes. To answer this question we present Theorem 3.1:

98 **Theorem 3.1.** *Given a dataset (X, Y) , and assuming small random initial network parameters, a*
 99 *small learning rate ϵ and that removable modes have smaller singular values than maintained modes,*
 100 *$G > 0$ (having multiple generations of learners) is a **necessary** condition for guaranteed removal of*
 101 *only the desired modes of variation.*

102 **Proof Sketch:** To prove this result we are required to show that the escaping time for the removable
 103 modes is greater than the hitting time of the maintained modes. We show that this is not true in general
 104 for one generation ($G=0$) using a contradicting example. Secondly, we show that after some number
 105 of generations G the removable mode escaping time is guaranteed to be larger than the maintained
 106 mode’s hitting time. The key step towards this is showing that the hitting time of the removable mode
 107 is larger than the hitting time of the maintained mode - a significantly easier condition than comparing
 108 hitting time and escaping time but which is enough for IL to be applicable. Once IL is applicable
 109 then as $G \rightarrow \infty$ the removable mode escaping time will become larger than the maintained mode
 110 hitting time, proving the theorem. The full proof is shown in Section B.

111 By assuming that the modes we aim to maintain are learned faster than removable modes we have
 112 imparted a preference in which modes of variation are learned and their relative ordering. This
 113 may appear to be a strong assumption, however, to maintain the slower learning modes would be
 114 a fundamental disagreement with IL as an algorithm. IL assumes that the fastest learning modes are
 115 systematic and provide the best generalization. Thus, for Theorem 3.1 our assumption of the ordering
 116 of the SVs is no more strict than assuming that IL is a valid algorithm for the dataset. We analyse
 117 the validity and limitations of the IL assumption that the quicker modes are systematic in Section 3.3.

118 3.2 The Requirement of Depth

119 A similar derivation of the IL learning dynamics can be done for a shallow network (no hidden layer).
 120 In this case the singular values of the model’s mapping follow the trajectory:

$$\pi_\alpha^G(t) = \pi_\alpha^{G-1} (1 - \exp(-\delta_\alpha t/\tau)) + \pi_0 \exp(-\delta_\alpha t/\tau), \quad (3)$$

121 such that the time course depends on the singular values of the input covariance matrix, δ_α (Saxe
 122 et al., 2019). While deep networks show stage-like transitions which allow for one mode to be
 123 learned while another remains near 0; in shallow networks the modes show an exponential approach
 124 to their asymptote and all modes are learned at once. Thus, there will never be an opportunity for
 125 IL to remove a mode without also losing information on modes which we aim to maintain. See the
 126 simulated runs in Figure 2 for a visual depiction of this fact. This mean that for IL to be an effective
 127 procedure for the refinement of language, depth is required in the network architecture.

128 3.3 IL Uncovers Systematic Modules

129 To establish whether IL has a benefit for systematicity we must formalize a space of datasets
 130 which display the inductive biases of learning. We build on prior work which aimed to formalize
 131 systematicity (Jarvis et al., 2023) and provides such a space of datasets. To be applicable to the
 132 linguistic background of IL we phrase the space of datasets in terms of the mapping from logical
 133 forms (our internal, potentially semantic, representation of the world) to phonetic forms (words or
 134 sentences) commonly discussed in linguistics (Brighton & Kirby, 2006). The fundamental aspect
 135 of the analysis, however, remains the same: we use a space of datasets parametrized by the degree
 136 of compositional and non-compositional structure. We then use the closed-form SVD for all datasets
 137 in the space (written in terms of the dataset parameters) to establish how dataset structure affects
 138 the inductive bias of the neural network learning dynamics.

139 To formalize this setting Jarvis et al. (2023) define a parametric space of datasets with input and output
 140 matrices $X = [X_\Omega \ X_\Gamma]^T$ and $Y = [Y_\Omega \ Y_\Gamma]^T$ respectively, where $n_x, n_y, k_x, k_y, r \in \mathbb{Z}^+$ are the pa-
 141 rameters that define a specific dataset. The *compositional input feature* matrix $X_\Omega \in \{-1, 1\}^{2^{n_x} \times n_x}$
 142 consists of all binary patterns with n_x bits. Here n_x is a key parameter determining the number
 143 of bits in the compositional input structure. Overall, the dataset contains 2^{n_x} examples. The
 144 *non-compositional input feature* matrix $X_\Gamma = [rI_1 \dots rI_{k_x}]$ consists of k_x scaled identity matrices,
 145 $I_i \in \{0, 1\}^{2^{n_x} \times 2^{n_x}}$ giving each datapoint a unique non-compositional identifying feature. Similar
 146 matrices are then defined for the output space. As described above, the network’s total input-output
 147 mapping at all times in training is a function of the singular value decomposition of the dataset
 148 statistics. For all datasets in the space there are three distinct input-output covariance (Σ^{yx}) singular
 149 values λ_1, λ_2 and λ_3 :

$$\lambda_1 = \left(\frac{(k_x r^2 + 2^{n_x})(k_y r^2 + 2^{n_x})}{2^{2n_x}} \right)^{\frac{1}{2}} \quad \lambda_2 = \left(\frac{(k_x r^2 + 2^{n_x})(k_y r^2)}{2^{2n_x}} \right)^{\frac{1}{2}} \quad \lambda_3 = \left(\frac{k_x k_y r^4}{2^{2n_x}} \right)^{\frac{1}{2}}$$

150 Note that the singular values are written in terms of the five dataset parameters, which allows for an
 151 analysis of how dataset structure influences the network training dynamics (Jarvis et al., 2023) and
 152 inductive bias of IL. Substituting these expressions into the dynamics equations we obtain equations
 153 for the networks mapping and full learning trajectories across all generations, at all times in training
 154 and for all datasets in the space. Appendix D depicts simulation to verify our theoretical results and
 155 we see exact agreement between the predicted and simulated dynamics. From this we can prove that
 156 for all datasets in the space the singular values will begin to learn in order. This means that IL can
 157 extract earlier singular values and remove later one. Since for the space of datasets the compositional
 158 input and output structure is only connected to the first mode π_1 (see Figure 1) IL is able to extract
 159 compositional, low-rank structure from the dataset. By the definition of systematicity in Jarvis et al.
 160 (2023), the reliance on low-rank structure here means that this module is generalizing systematically.
 161 This is a promising result and we hope that this will motivate more uses of IL in practical ML
 162 algorithms, for example like in the recent work of Ren et al. (2024).

163 **Observation 3.2.** *For all points in the space of datasets: $n_x, n_y, k_x, k_y, r \in \mathbb{Z}^+$ the input-output*
 164 *covariance matrix Σ^{yx} singular values will be ordered as: $\lambda_1 > \lambda_2 > \lambda_3$.*

165 **Proof Sketch:** The proof of this observation shows that there is no configuration of dataset parameters
 166 such that Equation λ_1 is not the largest value and λ_3 is not the smallest. We begin by assuming
 167 that this ordering holds and then simplify the expressions until we arrive at the requirement that
 168 $n_x, n_y, k_x, k_y, r \in \mathbb{Z}^+$ which is true by definition. The full proof is shown in Section C.

169 **References**

- 170 Jacob Andreas, Marcus Rohrbach, Trevor Darrell, and Dan Klein. Neural module networks. In
171 *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 39–48, 2016.
- 172 S. Arora, N. Cohen, and E. Hazan. On the optimization of deep networks: Implicit acceleration
173 by overparameterization. *35th International Conference on Machine Learning, ICML 2018*, 1:
174 372–389, 2018. arXiv: 1802.06509 ISBN: 9781510867963.
- 175 Dzmitry Bahdanau, Shikhar Murty, Michael Noukhovitch, Thien Huu Nguyen, Harm de Vries,
176 and Aaron Courville. Systematic generalization: What is required and can it be learned? In
177 *International Conference on Learning Representations*, 2019.
- 178 P. Baldi and K. Hornik. Neural networks and principal component analysis: Learning from examples
179 without local minima. *Neural Networks*, 2(1):53–58, January 1989. ISSN 08936080. doi:
180 10.1016/0893-6080(89)90014-2. URL [http://linkinghub.elsevier.com/retrieve/pii/
181 0893608089900142](http://linkinghub.elsevier.com/retrieve/pii/0893608089900142).
- 182 James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal
183 Maclaurin, George Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, and
184 Qiao Zhang. JAX: composable transformations of Python+NumPy programs, 2018. URL
185 <http://github.com/google/jax>.
- 186 Henry Brighton and Simon Kirby. Understanding linguistic evolution by visualizing the emergence
187 of topographic mappings. *Artificial life*, 12(2):229–242, 2006.
- 188 Jerry A Fodor and Zenon W Pylyshyn. Connectionism and cognitive architecture: A critical analysis.
189 *Cognition*, 28(1-2):3–71, 1988.
- 190 K. Fukumizu. Effect of Batch Learning In Multilayer Neural Networks. In *Proceedings of the 5th
191 International Conference on Neural Information Processing*, pp. 67–70, 1998.
- 192 Mario Geiger, Stefano Spigler, Arthur Jacot, and Matthieu Wyart. Disentangling feature and lazy
193 training in deep neural networks. *Journal of Statistical Mechanics: Theory and Experiment*, 2020
194 (11):113301, 2020.
- 195 Robert F Hadley. Connectionism, explicit rules, and symbolic manipulation. *Minds and machines*, 3
196 (2):183–200, 1993.
- 197 Charles F Hockett and Charles D Hockett. The origin of speech. *Scientific American*, 203(3):88–97,
198 1960.
- 199 Ronghang Hu, Jacob Andreas, Marcus Rohrbach, Trevor Darrell, and Kate Saenko. Learning to
200 reason: End-to-end module networks for visual question answering. In *Proceedings of the IEEE
201 International Conference on Computer Vision*, pp. 804–813, 2017.
- 202 Ronghang Hu, Jacob Andreas, Trevor Darrell, and Kate Saenko. Explainable neural computation via
203 stack neural module networks. In *Proceedings of the European conference on computer vision
204 (ECCV)*, pp. 53–69, 2018.
- 205 Minyoung Huh, Hossein Mobahi, Richard Zhang, Brian Cheung, Pulkit Agrawal, and Phillip Isola.
206 The low-rank simplicity bias in deep networks. *arXiv preprint arXiv:2103.10427*, 2021.
- 207 Dieuwke Hupkes, Verna Dankers, Mathijs Mul, and Elia Bruni. Compositionality decomposed: How
208 do neural networks generalise? *Journal of Artificial Intelligence Research*, 67:757–795, 2020.
- 209 Devon Jarvis, Richard Klein, Benjamin Rosman, and Andrew M Saxe. On the specialization of
210 neural modules. In *The Eleventh International Conference on Learning Representations*, 2023.
211 URL <https://openreview.net/forum?id=Fh97BDaR6I>.
- 212 Michael L Kalish, Thomas L Griffiths, and Stephan Lewandowsky. Iterated learning: Intergenerational
213 knowledge transmission reveals inductive biases. *Psychonomic Bulletin & Review*, 14(2):288–294,
214 2007.

- 215 Simon Kirby. Spontaneous evolution of linguistic structure-an iterated learning model of the emer-
 216 gence of regularity and irregularity. *IEEE Transactions on Evolutionary Computation*, 5(2):
 217 102–110, 2001.
- 218 Simon Kirby, Hannah Cornish, and Kenny Smith. Cumulative cultural evolution in the laboratory: An
 219 experimental approach to the origins of structure in human language. *Proceedings of the National
 220 Academy of Sciences*, 105(31):10681–10686, 2008.
- 221 Simon Kirby, Monica Tamariz, Hannah Cornish, and Kenny Smith. Compression and communication
 222 in the cultural evolution of linguistic structure. *Cognition*, 141:87–102, 2015.
- 223 Brenden Lake and Marco Baroni. Generalization without systematicity: On the compositional skills
 224 of sequence-to-sequence recurrent networks. In *Proceedings of the 35th International Conference
 225 on Machine Learning*, pp. 4487–4499. International Machine Learning Society (IMLS), 2018.
- 226 Brenden M Lake, Tomer D Ullman, Joshua B Tenenbaum, and Samuel J Gershman. Building
 227 machines that learn and think like people. *Behavioral and brain sciences*, 40, 2017.
- 228 Brenden M Lake, Tal Linzen, and Marco Baroni. Human few-shot learning of compositional
 229 instructions. In *Proceedings of the 41st Annual Conference of the Cognitive Science Society*, 2019.
- 230 A.K. Lampinen and S. Ganguli. An analytic theory of generalization dynamics and transfer learning
 231 in deep linear networks. In T. Sainath (ed.), *International Conference on Learning Representations*,
 232 2019. ISBN 0311-5518. doi: 10.1080/03115519808619195. URL [http://arxiv.org/abs/
 233 1809.10374](http://arxiv.org/abs/1809.10374). arXiv: 1809.10374.
- 234 Angeliki Lazaridou, Karl Moritz Hermann, Karl Tuyls, and Stephen Clark. Emergence of linguistic
 235 communication from referential games with symbolic and pixel input. In *International Conference
 236 on Learning Representations*, 2018.
- 237 Yuchen Lu, Soumye Singhal, Florian Strub, Aaron Courville, and Olivier Pietquin. Countering
 238 language drift with seeded iterated learning. In *International Conference on Machine Learning*,
 239 pp. 6437–6447. PMLR, 2020a.
- 240 Yuchen Lu, Soumye Singhal, Florian Strub, Olivier Pietquin, and Aaron Courville. Supervised
 241 seeded iterated learning for interactive language learning. In *Proceedings of the 2020 Conference
 242 on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 3962–3970, 2020b.
- 243 Lars Niklasson and Noel Sharkey. *Systematicity and generalisation in connectionist compositional
 244 representations*. Citeseer, 1992.
- 245 Steven Phillips and Janet Wiles. Exponential generalizations from a polynomial number of examples
 246 in a combinatorial domain. In *Proceedings of 1993 International Conference on Neural Networks
 247 (IJCNN-93-Nagoya, Japan)*, volume 1, pp. 505–508. IEEE, 1993.
- 248 Steven Andrew Phillips. *Connectionism and the problem of systematicity*. PhD thesis, University of
 249 Queensland, 1995.
- 250 Jordan B Pollack. Recursive distributed representations. *Artificial Intelligence*, 46(1-2):77–105,
 251 1990.
- 252 Yi Ren, Shangmin Guo, Matthieu Labeau, Shay B Cohen, and Simon Kirby. Compositional lan-
 253 guages emerge in a neural iterated learning model. In *International Conference on Learning
 254 Representations*, 2019.
- 255 Yi Ren, Samuel Lavoie, Michael Galkin, Danica J Sutherland, and Aaron C Courville. Improving
 256 compositional generalization using iterated learning and simplicial embeddings. *Advances in
 257 Neural Information Processing Systems*, 36, 2024.
- 258 A.M. Saxe, J.L. McClelland, and S. Ganguli. Exact solutions to the nonlinear dynamics of learning
 259 in deep linear neural networks. In Y. Bengio and Y. LeCun (eds.), *International Conference on
 260 Learning Representations*, Banff, Canada, 2014. Oral presentation. arXiv: 1312.6120v3.

- 261 Andrew M Saxe, James L McClelland, and Surya Ganguli. A mathematical theory of semantic
262 development in deep neural networks. *Proceedings of the National Academy of Sciences*, 116(23):
263 11537–11546, 2019.
- 264 Ankit Vani, Max Schwarzer, Yuchen Lu, Eeshan Dhekane, and Aaron Courville. Iterated learning for
265 emergent systematicity in vqa. In *International Conference on Learning Representations*, 2021.

266 **A Linear Dynamics Review**

267 While deep linear networks can only represent linear input-output mappings, the dynamics of learning
 268 change dramatically with the introduction of one or more hidden layers (Fukumizu, 1998; Saxe et al.,
 269 2014, 2019; Arora et al., 2018; Lampinen & Ganguli, 2019), and the learning problem becomes
 270 non-convex (Baldi & Hornik, 1989). They therefore serve as a tractable model of the influence of
 271 depth specifically on learning dynamics, which prior work has shown to impart a low-rank inductive
 272 bias on the linear mapping (Huh et al., 2021).

273 Consider a single hidden layer network computing output $\hat{y} = W^2 W^1 x$ in response to an input x ,
 274 trained to minimize the quadratic loss $L(W^1, W^2) = \frac{1}{2} \|Y - W^2 W^1 X\|_2^2$ using full batch gradient
 275 descent. This gives the learning rules for each layer as $E[\Delta W^1] = \epsilon W^{2T} (Y - W^2 W^1 X) X^T$ and
 276 $E[\Delta W^2] = \epsilon (Y - W^2 W^1 X) (W^1 X)^T$. By using a small learning rate ϵ and taking the continuous
 277 time limit, the mean change in weights is given by $\tau \frac{d}{dt} W^1 = W^{2T} (\Sigma^{yx} - W^2 W^1 \Sigma^x)$ and $\tau \frac{d}{dt} W^2 =$
 278 $(\Sigma^{yx} - W^2 W^1 \Sigma^x) W^{1T}$ where $\Sigma^x = E[XX^T]$ is the input correlation matrix, $\Sigma^{yx} = E[YX^T]$ is
 279 the input-output correlation matrix and $\tau = \frac{1}{P\epsilon}$ is the learning time constant for P inputs. Here, t
 280 measures units of learning epochs. It is helpful to note that since we are using a small learning rate
 281 the full batch gradient descent and stochastic gradient descent dynamics will be the same. Saxe et
 282 al. (2019) Saxe et al. (2019) has shown that the learning dynamics depend on the singular value
 283 decomposition of

$$\begin{aligned} \Sigma^x &= V D V^T = \sum_{\alpha=1}^{|X|} \delta_\alpha u^\alpha v^{\alpha T}; \\ \Sigma^{yx} &= U S V^T = \sum_{\alpha=1}^{\min(|X|, |Y|)} \lambda_\alpha u^\alpha v^{\alpha T} \end{aligned} \quad (4)$$

284 where U and V are orthogonal matrices of singular vectors and S, D are diagonal matrices of singular
 285 values/eigenvalues. To solve for the dynamics we require that the right singular vectors V of Σ^{yx} are
 286 also the singular vectors of Σ^x . We also assume that $n_h > \min(|X|, |Y|)$ where n_h is the number of
 287 hidden neurons. If this is not the case then the model will only learn the top n_h singular values of
 288 the input-output mapping (Saxe et al., 2014). Given the SVDs of the two correlation matrices the
 289 learning dynamics can be described explicitly as

$$W^2(t) W^1(t) = U A(t) V^T = \sum_{\alpha=1}^{2^{n_x}} \pi_\alpha(t) u^\alpha v^{\alpha T} \quad (5)$$

290 where $A(t)$ is the effective singular value matrix of the network's mapping.

291 **B Derivation of Escaping and Hitting Time**

292 We will begin by using the definition of escaping time and substituting the mode dynamics (Equation
 293 4) into this expression to obtain an expression for t :

$$\begin{aligned}
 \pi_\alpha^G &= \rho \\
 \frac{\lambda_\alpha^G / \delta_\alpha}{1 - \left(1 - \frac{\lambda_\alpha^G}{\delta_\alpha \pi_0}\right) \exp\left(\frac{-2\lambda_\alpha^G}{\tau} t\right)} &= \rho \\
 \frac{\lambda_\alpha^G / \delta_\alpha}{\rho} &= 1 - \left(1 - \frac{\lambda_\alpha^G}{\delta_\alpha \pi_0}\right) \exp\left(\frac{-2\lambda_\alpha^G}{\tau} t\right) \\
 \frac{\rho - \lambda_\alpha^G / \delta_\alpha}{\rho} &= \left(1 - \frac{\lambda_\alpha^G}{\delta_\alpha \pi_0}\right) \exp\left(\frac{-2\lambda_\alpha^G}{\tau} t\right) \\
 \frac{\rho - \lambda_\alpha^G / \delta_\alpha}{\rho \left(1 - \frac{\lambda_\alpha^G}{\delta_\alpha \pi_0}\right)} &= \exp\left(\frac{-2\lambda_\alpha^G}{\tau} t\right) \\
 \log\left(\frac{\rho - \lambda_\alpha^G / \delta_\alpha}{\rho \left(1 - \frac{\lambda_\alpha^G}{\delta_\alpha \pi_0}\right)}\right) &= \frac{-2\lambda_\alpha^G}{\tau} t \\
 \frac{-\tau}{2\lambda_\alpha^G} \log\left(\frac{\rho - \lambda_\alpha^G / \delta_\alpha}{\rho \left(1 - \frac{\lambda_\alpha^G}{\delta_\alpha \pi_0}\right)}\right) &= t \\
 \frac{-\tau}{2\lambda_\alpha^G} \log\left(\frac{\rho \left(1 - \frac{\lambda_\alpha^G}{\delta_\alpha \rho}\right)}{\rho \left(1 - \frac{\lambda_\alpha^G}{\delta_\alpha \pi_0}\right)}\right) &= t \\
 \frac{\tau}{2\lambda_\alpha^G} \log\left(\frac{1 - \frac{\lambda_\alpha^G}{\delta_\alpha \pi_0}}{1 - \frac{\lambda_\alpha^G}{\delta_\alpha \rho}}\right) &= t \\
 \frac{\tau}{2\lambda_\alpha^G} \log\left(\frac{\frac{\lambda_\alpha^G}{\delta_\alpha \pi_0} - 1}{\frac{\lambda_\alpha^G}{\delta_\alpha \rho} - 1}\right) &= t \\
 \frac{\tau}{2\lambda_\alpha^G} \log\left(\frac{\frac{\lambda_\alpha^G - \delta_\alpha \pi_0}{\delta_\alpha \pi_0}}{\frac{\lambda_\alpha^G - \delta_\alpha \rho}{\delta_\alpha \rho}}\right) &= t \\
 \frac{\tau}{2\lambda_\alpha^G} \log\left(\frac{\lambda_\alpha^G \delta_\alpha \rho - \delta_\alpha^2 \pi_0 \rho}{\lambda_\alpha^G \delta_\alpha \pi_0 - \delta_\alpha^2 \pi_0 \rho}\right) &= t \\
 \frac{\tau}{2\lambda_\alpha^G} \log\left(\frac{\lambda_\alpha^G \rho - \delta_\alpha \pi_0 \rho}{\lambda_\alpha^G \pi_0 - \delta_\alpha \pi_0 \rho}\right) &= t
 \end{aligned}$$

294 The term inside of the log is always going to be greater than 1 as $\rho \geq \pi_0$. Thus the log is positive. In
 295 the extreme case where $\rho = \pi_0$ then the log evaluates to 0 as the internal fraction is 1, which makes
 296 sense as in this case the escaping time will be reached at initialization. We can perform a similar

297 computation for a modes hitting time.

$$\begin{aligned}
& (\lambda_\alpha^G/\delta_\alpha) - \pi_\alpha^G = \rho \\
& (\lambda_\alpha^G/\delta_\alpha) - \frac{\lambda_\alpha^G/\delta_\alpha}{1 - \left(1 - \frac{\lambda_\alpha^G}{\delta_\alpha\pi_0}\right) \exp\left(\frac{-2\lambda_\alpha^G}{\tau}t\right)} = \rho \\
& \frac{\lambda_\alpha^G/\delta_\alpha}{1 - \left(1 - \frac{\lambda_\alpha^G}{\delta_\alpha\pi_0}\right) \exp\left(\frac{-2\lambda_\alpha^G}{\tau}t\right)} = \rho - (\lambda_\alpha^G/\delta_\alpha) \\
& \frac{\lambda_\alpha^G/\delta_\alpha}{\rho - (\lambda_\alpha^G/\delta_\alpha)} = 1 - \left(1 - \frac{\lambda_\alpha^G}{\delta_\alpha\pi_0}\right) \exp\left(\frac{-2\lambda_\alpha^G}{\tau}t\right) \\
& \frac{\lambda_\alpha^G/\delta_\alpha}{\rho - (\lambda_\alpha^G/\delta_\alpha)} - \frac{\rho - (\lambda_\alpha^G/\delta_\alpha)}{\rho - (\lambda_\alpha^G/\delta_\alpha)} = - \left(1 - \frac{\lambda_\alpha^G}{\delta_\alpha\pi_0}\right) \exp\left(\frac{-2\lambda_\alpha^G}{\tau}t\right) \\
& \frac{\rho}{\rho - (\lambda_\alpha^G/\delta_\alpha)} = \left(1 - \frac{\lambda_\alpha^G}{\delta_\alpha\pi_0}\right) \exp\left(\frac{-2\lambda_\alpha^G}{\tau}t\right) \\
& \frac{\rho}{(\rho - (\lambda_\alpha^G/\delta_\alpha)) \left(1 - \frac{\lambda_\alpha^G}{\delta_\alpha\pi_0}\right)} = \exp\left(\frac{-2\lambda_\alpha^G}{\tau}t\right) \\
& \log\left(\frac{\rho}{(\rho - (\lambda_\alpha^G/\delta_\alpha)) \left(1 - \frac{\lambda_\alpha^G}{\delta_\alpha\pi_0}\right)}\right) = \frac{-2\lambda_\alpha^G}{\tau}t \\
& \frac{-\tau}{2\lambda_\alpha^G} \log\left(\frac{\rho}{((\lambda_\alpha^G/\delta_\alpha) - \rho) \left(\frac{\lambda_\alpha^G}{\delta_\alpha\pi_0} - 1\right)}\right) = t \\
& \frac{\tau}{2\lambda_\alpha^G} \log\left(\frac{((\lambda_\alpha^G/\delta_\alpha) - \rho) \left(\frac{\lambda_\alpha^G}{\delta_\alpha\pi_0} - 1\right)}{\rho}\right) = t \\
& \frac{\tau}{2\lambda_\alpha^G} \log\left(\frac{\lambda_\alpha^{G^2}}{\delta_\alpha^2\pi_0\rho} - \frac{\lambda_\alpha^G}{\delta_\alpha\rho} - \frac{\lambda_\alpha^G}{\delta_\alpha\pi_0} + 1\right) = t \\
& \frac{\tau}{2\lambda_\alpha^G} \log\left(\frac{\lambda_\alpha^{G^2} - \delta_\alpha\lambda_\alpha^G\pi_0 - \delta_\alpha\lambda_\alpha^G\rho + \delta_\alpha^2\pi_0\rho}{\delta_\alpha^2\pi_0\rho}\right) = t
\end{aligned}$$

298 Thus we can summarize the escaping time and hitting time for mode α as follows:

$$\hat{t}_\alpha = \frac{\tau}{2\lambda_\alpha^G} \log\left(\frac{\lambda_\alpha^G\rho - \delta_\alpha\pi_0\rho}{\lambda_\alpha^G\pi_0 - \delta_\alpha\pi_0\rho}\right); t_\alpha^* = \frac{\tau}{2\lambda_\alpha^G} \log\left(\frac{\lambda_\alpha^{G^2} - \delta_\alpha\lambda_\alpha^G\pi_0 - \delta_\alpha\lambda_\alpha^G\rho + \delta_\alpha^2\pi_0\rho}{\delta_\alpha^2\pi_0\rho}\right)$$

299 Proof of Theorem 3.1

300 **Theorem 3.1.** Given a dataset (X, Y) , and assuming small random initial network parameters, a
301 small learning rate ϵ and that removable modes have smaller singular values than maintained modes,
302 $G > 0$ (having multiple generations of learners) is a **necessary** condition for guaranteed removal of
303 only the desired modes of variation.

304 We note that for now we make no assumptions on the dataset or singular values aside from their
305 usual ordering. In subsequent sections we will make reference to the particular space of datasets.
306 We now aim to show that there will always be a $G \geq 0$ for which the removable mode's hitting
307 time will be larger than the maintained modes convergence time. To begin we need to show that
308 the convergence time of the maintained mode is larger than the convergence time of the removable
309 mode for all generations. This means that it is possible to find a time where the maintained mode has
310 converged but the removable mode has not and IL is a valid algorithm for removing the removable
311 mode. We may keep G general here as the generation does not change the proof since a removable

312 mode can only decrease in size. Thus, we begin by determining under what conditions the maintained
 313 mode will have an earlier convergence time than the removable modes:

$$\frac{\tau}{2\lambda_{\alpha-1}^G} \log \left(\frac{\lambda_{\alpha-1}^{G^2} - \delta_{\alpha-1}\lambda_{\alpha-1}^G\pi_0 - \delta_{\alpha-1}\lambda_{\alpha-1}^G\rho + \delta_{\alpha-1}^2\pi_0\rho}{\delta_{\alpha-1}^2\pi_0\rho} \right) > \frac{\tau}{2\lambda_{\alpha}^G} \log \left(\frac{\lambda_{\alpha}^{G^2} - \delta_{\alpha}\lambda_{\alpha}^G\pi_0 - \delta_{\alpha}\lambda_{\alpha}^G\rho + \delta_{\alpha}^2\pi_0\rho}{\delta_{\alpha}^2\pi_0\rho} \right)$$

314 We next substitute in the fact that $\pi_0 \rightarrow 0$ and $\rho \rightarrow 0$. However, we need to consider first that $\lambda_{\alpha-1}^G$
 315 may be small (especially when IL is being applied). The above inequality will only hold when the
 316 log on the left of the inequality is positive. Thus, we first consider when the argument to the log is
 317 greater than 1:

$$\begin{aligned} \frac{\lambda_{\alpha-1}^{G^2} - \delta_{\alpha-1}\lambda_{\alpha-1}^G\pi_0 - \delta_{\alpha-1}\lambda_{\alpha-1}^G\rho + \delta_{\alpha-1}^2\pi_0\rho}{\delta_{\alpha-1}^2\pi_0\rho} &> 1 \\ \frac{\lambda_{\alpha-1}^{G^2} - \delta_{\alpha-1}\lambda_{\alpha-1}^G\pi_0 - \delta_{\alpha-1}\lambda_{\alpha-1}^G\rho}{\delta_{\alpha-1}^2\pi_0\rho} + 1 &> 1 \\ \frac{\lambda_{\alpha-1}^{G^2} - \delta_{\alpha-1}\lambda_{\alpha-1}^G\pi_0 - \delta_{\alpha-1}\lambda_{\alpha-1}^G\rho}{\delta_{\alpha-1}^2\pi_0\rho} &> 0 \\ \lambda_{\alpha-1}^G - \delta_{\alpha-1}\pi_0 - \delta_{\alpha-1}\rho &> 0 \\ \frac{\lambda_{\alpha-1}^G}{\delta_{\alpha-1}} &> \pi_0 + \rho \end{aligned}$$

318 Thus, for the log to be positive the final value of the mode must be larger than $\pi_0 + \rho$. This is a very
 319 easy constraint to meet and any mode with a final value less than $\pi_0 + \rho$ will have no real bearing on
 320 the output language. Additionally, in the limits of $\pi_0 \rightarrow 0$ and $\rho \rightarrow 0$ this just means that the final
 321 mode must be greater than 0 which is true by definition. We are free to substitute in $\pi_0 \rightarrow 0$ and
 322 $\rho \rightarrow 0$ to the original expression and it simplifies to the following where $c = \frac{1}{\pi_0\rho}$.

$$\begin{aligned} \frac{1}{2\lambda_{\alpha-1}^G} \log \left(c \frac{\lambda_{\alpha-1}^{G^2}}{\delta_{\alpha-1}^2} \right) &> \frac{1}{2\lambda_{\alpha}^G} \log \left(c \frac{\lambda_{\alpha}^{G^2}}{\delta_{\alpha}^2} \right) \\ \frac{1}{\lambda_{\alpha-1}^G} \log \left(\sqrt{c} \frac{\lambda_{\alpha-1}^G}{\delta_{\alpha-1}} \right) &> \frac{1}{\lambda_{\alpha}^G} \log \left(\sqrt{c} \frac{\lambda_{\alpha}^G}{\delta_{\alpha}} \right) \end{aligned}$$

323 What this demonstrates is that, with small initial mode values and a high precision, the convergence
 324 time is inversely related to the size of the input-output covariance of the mode and proportional to the
 325 log of the final value of the mode itself. In essence, how long it takes learning to converge depends
 326 on how quickly learning happens and how much there is to learn. Noting that we can set c inside
 327 of the log and know that it is a very large value ($c \rightarrow \infty$) we can put the expression in a regime
 328 where the log derivative is near 0 and we can treat the two log expressions as constant and equal. We
 329 note that this is only valid where the log is positive which we have demonstrated is the case. How
 330 large we need to set c in practice depends entirely on the relative scale of $\lambda_{\alpha-1}^G/\delta_{\alpha-1}$ and $\lambda_{\alpha}^G/\delta_{\alpha}$.
 331 If $\lambda_{\alpha-1}^G/\delta_{\alpha-1}$ and $\lambda_{\alpha}^G/\delta_{\alpha}$ are of roughly the same magnitude then it may not even be necessary to
 332 set a large c for the inequality to hold, but this then depends on the relative scale of $\lambda_{\alpha-1}^G$ and λ_{α}^G
 333 in isolation too. What we are accounting for here is the case where a slower mode also has a much
 334 smaller final value, in which case it may still converge as quickly as the faster learning mode. The
 335 way to ensure these modes maintain their ordering is to set a smaller initial parameters and a smaller
 336 precision. This works because the dynamics have the sigmoidal training curve where overcoming the
 337 initial saddle point takes long and then there is a stage-like transition once the mode begins being
 338 learned. This simplifies the expression further to:

$$\frac{1}{\lambda_{\alpha-1}^G} \gtrsim \frac{1}{\lambda_{\alpha}^G}$$

339 This is true by definition and so it is appropriate to apply IL for all generations. We will now consider
 340 the relationship between the removable modes hitting time and the maintained modes convergence

341 time.

$$\hat{t}_{\alpha-1} > t_{\alpha}^*$$

$$\frac{\tau}{2\lambda_{\alpha-1}^G} \log \left(\frac{\lambda_{\alpha-1}^G \rho - \delta_{\alpha-1} \pi_0 \rho}{\lambda_{\alpha-1}^G \pi_0 - \delta_{\alpha-1} \pi_0 \rho} \right) > \frac{\tau}{2\lambda_{\alpha}^G} \log \left(\frac{\lambda_{\alpha}^{G^2} - \delta_{\alpha} \lambda_{\alpha}^G \pi_0 - \delta_{\alpha} \lambda_{\alpha}^G \rho + \delta_{\alpha}^2 \pi_0 \rho}{\delta_{\alpha}^2 \pi_0 \rho} \right)$$

342 Once again, applying the usual limits of $\pi_0 \rightarrow 0$ and $\rho \rightarrow 0$ simplifies the expression to:

$$\begin{aligned} \frac{1}{\lambda_{\alpha-1}^G} \log \left(\frac{\lambda_{\alpha-1}^G \rho}{\lambda_{\alpha-1}^G \pi_0} \right) &> \frac{1}{\lambda_{\alpha}^G} \log \left(c \frac{\lambda_{\alpha}^{G^2}}{\delta_{\alpha}^2} \right) \\ \frac{1}{\lambda_{\alpha-1}^G} \log \left(\frac{\rho}{\pi_0} \right) &> \frac{2}{\lambda_{\alpha}^G} \log \left(\sqrt{c} \frac{\lambda_{\alpha}^G}{\delta_{\alpha}} \right) \\ \frac{1}{\lambda_{\alpha-1}^G} \log \left(\frac{\rho}{\pi_0} \right) &> \frac{2}{\lambda_{\alpha}^G} \log \left(\sqrt{c} \frac{\lambda_{\alpha}^G}{\delta_{\alpha}} \right) \\ \frac{1}{\lambda_{\alpha-1}^G} \log (\rho^2 c) &> \frac{2}{\lambda_{\alpha}^G} \log \left(\sqrt{c} \frac{\lambda_{\alpha}^G}{\delta_{\alpha}} \right) \\ \frac{2}{\lambda_{\alpha-1}^G} \log (\rho \sqrt{c}) &> \frac{2}{\lambda_{\alpha}^G} \log \left(\sqrt{c} \frac{\lambda_{\alpha}^G}{\delta_{\alpha}} \right) \\ \frac{1}{\lambda_{\alpha-1}^G} \log (\rho \sqrt{c}) &> \frac{1}{\lambda_{\alpha}^G} \log \left(\frac{\lambda_{\alpha}^G}{\delta_{\alpha}} \sqrt{c} \right) \\ \frac{1}{\lambda_{\alpha-1}^G} \log \left(\frac{\rho}{\sqrt{\rho \pi_0}} \right) &> \frac{1}{\lambda_{\alpha}^G} \log \left(\frac{\lambda_{\alpha}^G}{\delta_{\alpha}} \frac{1}{\sqrt{\rho \pi_0}} \right) \\ \frac{1}{\lambda_{\alpha-1}^G} \log \left(\frac{\sqrt{\rho}}{\sqrt{\pi_0}} \right) &> \frac{1}{\lambda_{\alpha}^G} \log \left(\frac{\lambda_{\alpha}^G}{\delta_{\alpha}} \frac{1}{\sqrt{\rho \pi_0}} \right) \end{aligned}$$

343 This expression is not true in general. For example if we set: $\lambda_{\alpha-1}^G = 1$, $\lambda_{\alpha}^G = \sqrt{2}$ and $\delta_{\alpha} = 1$ then
344 we obtain:

$$\begin{aligned} \frac{1}{2} \log \left(\frac{\rho}{\pi_0} \right) &> \frac{1}{2\sqrt{2}} \log \left(\frac{2}{\rho \pi_0} \right) \\ \log \left(\frac{\rho}{\pi_0} \right) &> \frac{1}{\sqrt{2}} \log \left(\frac{2}{\rho \pi_0} \right) \end{aligned}$$

345 However $\rho \rightarrow 0$ and the left side of the expression is tending towards 0 while the right side is tending
346 towards ∞ . Thus, this is a contradiction. However, if we use the fact that $\lambda_{\alpha-1}^G \rightarrow 0$ as would be the
347 case when the IL algorithm is applied, then regardless of the other dataset statistics the expression
348 simplifies to:

$$\begin{aligned} \lim_{G \rightarrow \infty} \frac{1}{\lambda_{\alpha-1}^G} \log \left(\frac{\rho}{\pi_0} \right) &> \lim_{G \rightarrow \infty} \frac{2}{\lambda_{\alpha}^G} \log \left(\sqrt{c} \frac{\lambda_{\alpha}^G}{\delta_{\alpha}} \right) \\ \frac{1}{\lambda_{\alpha-1}^{\infty}} \log \left(\frac{\rho}{\pi_0} \right) &> \frac{2}{\lambda_{\alpha}^G} \log \left(\sqrt{c} \frac{\lambda_{\alpha}^G}{\delta_{\alpha}} \right) \\ &\infty > \frac{2}{\lambda_{\alpha}^G} \log \left(\sqrt{c} \frac{\lambda_{\alpha}^G}{\delta_{\alpha}} \right) \end{aligned}$$

349 This is true by definition. In practice all that is required is to find some G for which the expression
350 holds by a sufficient decrease in $\lambda_{\alpha-1}^G$. Thus, we have shown that it is always possible to reach a
351 point where a removable mode's hitting time is higher than the maintained modes convergence time.
352 This is not guaranteed to be the case for the first generation and may require multiple generations,
353 which have shown to always be a viable strategy. As a consequence it is always possible to remove a
354 removable mode while maintaining the maintained mode.

355 **C Proof of Observation 3.2**

356 **Observation 3.2.** For all points in the space of datasets: $n_x, n_y, k_x, k_y, r \in \mathbb{Z}^+$ the input-output
 357 covariance matrix Σ^{yx} singular values will be ordered as: $\lambda_1 > \lambda_2 > \lambda_3$.

358 Firstly we prove that $\lambda_1 > \lambda_2$:

$$\begin{aligned} \lambda_1 &> \lambda_2 \\ \left(\frac{(k_x r^2 + 2^{n_x})(k_y r^2 + 2^{n_x})}{2^{2n_x}} \right)^{\frac{1}{2}} &> \left(\frac{(k_x r^2 + 2^{n_x})(k_y r^2)}{2^{2n_x}} \right)^{\frac{1}{2}} \\ ((k_x r^2 + 2^{n_x})(k_y r^2 + 2^{n_x})) &> ((k_x r^2 + 2^{n_x})(k_y r^2)) \\ k_y r^2 + 2^{n_x} &> k_y r^2 \\ 2^{n_x} &> 0 \end{aligned}$$

359 $2^{n_x} > 0$ is true by definition since $n_x \in \mathbb{Z}^+$ and, thus, $\lambda_1 > \lambda_2$ for all points in our space of datasets.
 360 Now we prove that $\lambda_2 > \lambda_3$:

$$\begin{aligned} \lambda_2 &> \lambda_3 \\ \left(\frac{(k_x r^2 + 2^{n_x})(k_y r^2)}{2^{2n_x}} \right)^{\frac{1}{2}} &> \left(\frac{k_x k_y r^4}{2^{2n_x}} \right)^{\frac{1}{2}} \\ (k_x r^2 + 2^{n_x})(k_y r^2) &> k_x k_y r^4 \\ k_x k_y r^4 + 2^{n_x} k_y r^2 &> k_x k_y r^4 \\ 2^{n_x} k_y r^2 &> 0 \end{aligned}$$

361 $2^{n_x} k_y r^2 > 0$ is true by definition since $n_x, k_y, r \in \mathbb{Z}^+$ and, thus, $\lambda_2 > \lambda_3$ for all points in our space
 362 of datasets. Thus, using the transitivity of inequality: $\lambda_1 > \lambda_2 > \lambda_3$ for all points in the space of
 363 datasets.

364 **D Simulations of IL Dynamics**

365 To empirically verify our theoretical results we simulate the full training dynamics for deep and
 366 shallow linear networks trained using gradient descent on an instantiation from the space of datasets
 367 with parameters $n_x = 3, n_y = 2, k_x = 3, k_y = 1, r = 2$ (shown in Figure 2). While training,
 368 we compute the singular values of the network after each epoch. These simulations of the training
 369 dynamics for each unique singular value are then compared to the predicted dynamics. We also
 370 compute the Frobenius norms of portions of the network. These norms are functions of the singular
 371 values and summarize how entire portions of the input space connect to portions of the output space.
 372 Here we track how compositional/non-compositional inputs affect compositional/non-compositional
 373 outputs. The equations of these form Frobenius norms are shown below and were also first introduced
 374 in Jarvis et al. (2023). We see close agreement between the predicted and simulated trajectories¹. Note
 375 the requirement of depth and multiple generations of IL to effectively remove a mode of variation
 376 without also losing information on other modes. The difference between deep and shallow network
 377 training dynamics can be seen by comparing the shape of the learning trajectories between Figures
 378 2(a) and 2(c). Note how π_3 is removed with the deep network (Figure 2(a)) while π_1 and π_2 remain
 379 unchanged. In contrast, all modes are decreased with the shallow network (Figure 2(c)) but none are
 380 removed.

$$X_\Omega Y_\Omega^G\text{-Norm} = \left(\frac{2^{2n_x} n_y \pi_1^2(t)}{(k_x r^2 + 2^{n_x})(k_y r^2 + 2^{n_x})} \right)^{\frac{1}{2}} \quad (6)$$

$$X_\Gamma Y_\Gamma^G\text{-Norm} = \left(\frac{2^{n_x} n_y k_x r^2 \pi_1^2(t)}{(k_x r^2 + 2^{n_x})(k_y r^2 + 2^{n_x})} \right)^{\frac{1}{2}} \quad (7)$$

$$X_\Omega Y_\Gamma^G\text{-Norm} = \left(\frac{2^{n_x} k_y n_y r^2 \pi_1^2(t)}{(k_x r^2 + 2^{n_x})(k_y r^2 + 2^{n_x})} + \frac{2^{n_x} (n_x - n_y)}{k_x r^2 + 2^{n_x}} \pi_2^2(t) \right)^{\frac{1}{2}} \quad (8)$$

¹All experiments are run using the Jax library (Bradbury et al., 2018).

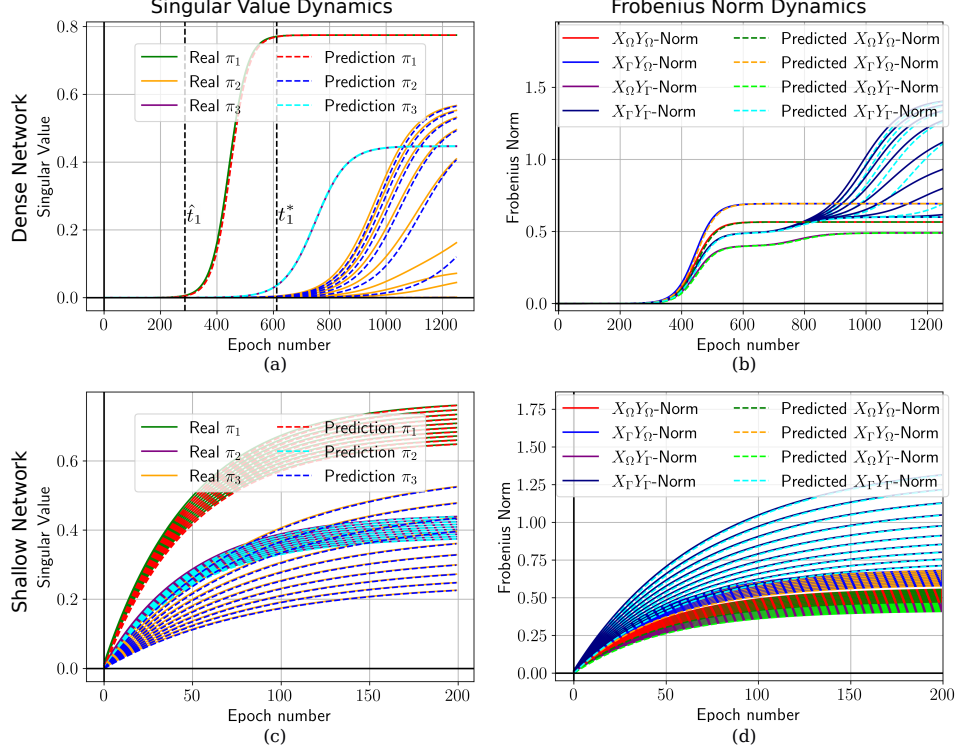


Figure 2: Analytical learning dynamics for deep (panels a-b) and shallow (panels c-d) linear networks. (a,c) Comparison of predicted and actual singular value trajectories over learning, for the three unique dataset singular values. \hat{t}_1 and t_1^* denote the escaping time and hitting time respectively for the first mode of variation with $\rho = 0.005$. (b,d) Comparison of predicted and actual Frobenius norms of the input-output mapping to/from compositional (X_Ω, Y_Ω) and non-compositional (X_Γ, Y_Γ) features. Deep networks show distinct stages of improvement over learning. However, at no point is a mapping learned which relies exclusively on compositional features or language. However, this setting depicts the progressive removal of the π_3 mode of variation over 10 generations. By the final generation of the dense network training the non-systematic norms exhibit two stage-like transitions corresponding to the learning of the two remaining modes of variation. The shallow network does not learn the modes in separate stages and so the removal of one distinct mode is impossible without simultaneously removing portions of all other modes. This demonstrates the theoretical observations from the dynamics of IL above. *Dataset Parameters:* $n_x = 3, n_y = 2, k_x = 3, k_y = 1, r = 2$.

383

$$X_\Gamma Y_\Gamma^G\text{-Norm} = \left(\frac{k_x k_y n_y r^4 \pi_1^2(t)}{(k_x r^2 + 2^{n_x})(k_y r^2 + 2^{n_x})} + \frac{(n_x - n_y) k_x r^2}{k_x r^2 + 2^{n_x}} \pi_2^2(t) + (2^{n_x} - n_x) \pi_3^2(t) \right)^{\frac{1}{2}} \quad (9)$$

384 E NeurIPS Paper Checklist

385 1. Claims

386 Question: Do the main claims made in the abstract and introduction accurately reflect the
387 paper's contributions and scope?

388 Answer: [Yes]

389 Justification: There are three main contributions in the abstract and each is given its own
390 section in our results. The introduction highlights the use of Saxe et al. (2019) and Jarvis
391 et al. (2023) in particular for informing or methods and this is exactly what we use.

392 Guidelines:

- 393 • The answer NA means that the abstract and introduction do not include the claims
394 made in the paper.
- 395 • The abstract and/or introduction should clearly state the claims made, including the
396 contributions made in the paper and important assumptions and limitations. A No or
397 NA answer to this question will not be perceived well by the reviewers.
- 398 • The claims made should match theoretical and experimental results, and reflect how
399 much the results can be expected to generalize to other settings.
- 400 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
401 are not attained by the paper.

402 2. Limitations

403 Question: Does the paper discuss the limitations of the work performed by the authors?

404 Answer: [Yes]

405 Justification: The limitations of the approach are clearly stated where appropriate. We also
406 review the linear dynamics and its assumptions in Section A.

407 Guidelines:

- 408 • The answer NA means that the paper has no limitation while the answer No means that
409 the paper has limitations, but those are not discussed in the paper.
- 410 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 411 • The paper should point out any strong assumptions and how robust the results are to
412 violations of these assumptions (e.g., independence assumptions, noiseless settings,
413 model well-specification, asymptotic approximations only holding locally). The authors
414 should reflect on how these assumptions might be violated in practice and what the
415 implications would be.
- 416 • The authors should reflect on the scope of the claims made, e.g., if the approach was
417 only tested on a few datasets or with a few runs. In general, empirical results often
418 depend on implicit assumptions, which should be articulated.
- 419 • The authors should reflect on the factors that influence the performance of the approach.
420 For example, a facial recognition algorithm may perform poorly when image resolution
421 is low or images are taken in low lighting. Or a speech-to-text system might not be
422 used reliably to provide closed captions for online lectures because it fails to handle
423 technical jargon.
- 424 • The authors should discuss the computational efficiency of the proposed algorithms
425 and how they scale with dataset size.
- 426 • If applicable, the authors should discuss possible limitations of their approach to
427 address problems of privacy and fairness.
- 428 • While the authors might fear that complete honesty about limitations might be used by
429 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
430 limitations that aren't acknowledged in the paper. The authors should use their best
431 judgment and recognize that individual actions in favor of transparency play an impor-
432 tant role in developing norms that preserve the integrity of the community. Reviewers
433 will be specifically instructed to not penalize honesty concerning limitations.

434 3. Theory Assumptions and Proofs

435 Question: For each theoretical result, does the paper provide the full set of assumptions and
436 a complete (and correct) proof?

437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491

Answer: [Yes]

Justification: We state the assumptions and provide a proof sketch in the main text. Full proves are given in the appendix.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental Result Reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: We describe the datasets, network architectures and training algorithms for all experiments in the work.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not. Finally, we release the code for all experiments (including bash scripts to make this easy) with a requirements file for reproducibility.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [No]

Justification: This is a preprint, and so we omit code at this point. However we guide a reader to code release for Jarvis et al. (2023) which provides code to imitate the setup of this work and the linear network dynamics in that setting.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental Setting/Details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: These were described in the description of the experimental design and background on the GDLN paradigm as this requires full-batch gradient descent. All other details were mentioned or are relatively easy to determine, such as choosing learning rate. In this case the learning rate has quite a broad range of valid choices as long as it is not set too large.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment Statistical Significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [NA]

Justification: All experiments in this work have deterministic results. The random initialization of the network parameters does not affect this as they are set small enough to be very low variance.

Guidelines:

- 544 • The answer NA means that the paper does not include experiments.
- 545 • The authors should answer "Yes" if the results are accompanied by error bars, confi-
- 546 dence intervals, or statistical significance tests, at least for the experiments that support
- 547 the main claims of the paper.
- 548 • The factors of variability that the error bars are capturing should be clearly stated (for
- 549 example, train/test split, initialization, random drawing of some parameter, or overall
- 550 run with given experimental conditions).
- 551 • The method for calculating the error bars should be explained (closed form formula,
- 552 call to a library function, bootstrap, etc.)
- 553 • The assumptions made should be given (e.g., Normally distributed errors).
- 554 • It should be clear whether the error bar is the standard deviation or the standard error
- 555 of the mean.
- 556 • It is OK to report 1-sigma error bars, but one should state it. The authors should
- 557 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
- 558 of Normality of errors is not verified.
- 559 • For asymmetric distributions, the authors should be careful not to show in tables or
- 560 figures symmetric error bars that would yield results that are out of range (e.g. negative
- 561 error rates).
- 562 • If error bars are reported in tables or plots, The authors should explain in the text how
- 563 they were calculated and reference the corresponding figures or tables in the text.

564 8. Experiments Compute Resources

565 Question: For each experiment, does the paper provide sufficient information on the com-
 566 puter resources (type of compute workers, memory, time of execution) needed to reproduce
 567 the experiments?

568 Answer: [Yes]

569 Justification: The networks we use here are very small and train on the order of seconds
 570 on a single Nvidia 1080. The necessary compute is intuitive based off of the architecture
 571 descriptions.

572 Guidelines:

- 573 • The answer NA means that the paper does not include experiments.
- 574 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
- 575 or cloud provider, including relevant memory and storage.
- 576 • The paper should provide the amount of compute required for each of the individual
- 577 experimental runs as well as estimate the total compute.
- 578 • The paper should disclose whether the full research project required more compute
- 579 than the experiments reported in the paper (e.g., preliminary or failed experiments that
- 580 didn't make it into the paper).

581 9. Code Of Ethics

582 Question: Does the research conducted in the paper conform, in every respect, with the
 583 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

584 Answer: [Yes]

585 Justification: This is a theory work and so we do not perceive of any potential ethics
 586 concerns.

587 Guidelines:

- 588 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 589 • If the authors answer No, they should explain the special circumstances that require a
- 590 deviation from the Code of Ethics.
- 591 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
- 592 eration due to laws or regulations in their jurisdiction).

593 10. Broader Impacts

594 Question: Does the paper discuss both potential positive societal impacts and negative
 595 societal impacts of the work performed?

596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647

Answer: [NA]

Justification: Naturally we believe that our work will be of benefit to the research community which could have positive societal impact. However, our work is theoretical and to comment on any positive or negative societal impacts here would be speculative and potentially irresponsible. Since there are no primary or immediate societal impacts we omit any further discuss.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: None of our datasets or models appears to have a potential for misuse requiring safeguards.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

648 Justification: Citations are made to all works which directly influenced our own work or
649 approaches. We also cite the Jax library as the framework which we used for the experiments.
650 No other licensing issues arose.

651 Guidelines:

- 652 • The answer NA means that the paper does not use existing assets.
- 653 • The authors should cite the original paper that produced the code package or dataset.
- 654 • The authors should state which version of the asset is used and, if possible, include a
655 URL.
- 656 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 657 • For scraped data from a particular source (e.g., website), the copyright and terms of
658 service of that source should be provided.
- 659 • If assets are released, the license, copyright information, and terms of use in the
660 package should be provided. For popular datasets, `paperswithcode.com/datasets`
661 has curated licenses for some datasets. Their licensing guide can help determine the
662 license of a dataset.
- 663 • For existing datasets that are re-packaged, both the original license and the license of
664 the derived asset (if it has changed) should be provided.
- 665 • If this information is not available online, the authors are encouraged to reach out to
666 the asset's creators.

667 13. New Assets

668 Question: Are new assets introduced in the paper well documented and is the documentation
669 provided alongside the assets?

670 Answer: [NA]

671 Justification: No new asset have been made beyond the academic contribution.

672 Guidelines:

- 673 • The answer NA means that the paper does not release new assets.
- 674 • Researchers should communicate the details of the dataset/code/model as part of their
675 submissions via structured templates. This includes details about training, license,
676 limitations, etc.
- 677 • The paper should discuss whether and how consent was obtained from people whose
678 asset is used.
- 679 • At submission time, remember to anonymize your assets (if applicable). You can either
680 create an anonymized URL or include an anonymized zip file.

681 14. Crowdsourcing and Research with Human Subjects

682 Question: For crowdsourcing experiments and research with human subjects, does the paper
683 include the full text of instructions given to participants and screenshots, if applicable, as
684 well as details about compensation (if any)?

685 Answer: [NA]

686 Justification: No crowdsourcing or human subjects were used.

687 Guidelines:

- 688 • The answer NA means that the paper does not involve crowdsourcing nor research with
689 human subjects.
- 690 • Including this information in the supplemental material is fine, but if the main contribu-
691 tion of the paper involves human subjects, then as much detail as possible should be
692 included in the main paper.
- 693 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
694 or other labor should be paid at least the minimum wage in the country of the data
695 collector.

696 15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human 697 Subjects

698 Question: Does the paper describe potential risks incurred by study participants, whether
699 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
700 approvals (or an equivalent approval/review based on the requirements of your country or
701 institution) were obtained?

702 Answer: [NA]

703 Justification: No such concerns arose requiring review or approval.

704 Guidelines:

- 705 • The answer NA means that the paper does not involve crowdsourcing nor research with
706 human subjects.
- 707 • Depending on the country in which research is conducted, IRB approval (or equivalent)
708 may be required for any human subjects research. If you obtained IRB approval, you
709 should clearly state this in the paper.
- 710 • We recognize that the procedures for this may vary significantly between institutions
711 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
712 guidelines for their institution.
- 713 • For initial submissions, do not include any information that would break anonymity (if
714 applicable), such as the institution conducting the review.