Measuring and Controlling Solution Degeneracy across Task-Trained Recurrent Neural Networks

Anonymous Author(s)

Affiliation Address email

Abstract

Task-trained recurrent neural networks (RNNs) are widely used in neuroscience and machine learning to model dynamical computations. To gain mechanistic insight into how neural systems solve tasks, prior work often reverse-engineers individual trained networks. However, different RNNs trained on the same task and achieving similar performance can exhibit strikingly different internal solutions—a phenomenon known as solution degeneracy. Here, we develop a unified framework to systematically quantify and control solution degeneracy across three levels: behavior, neural dynamics, and weight space. We apply this framework to 3,400 RNNs trained on four neuroscience-relevant tasks while systematically varying task complexity, learning regime, network size, and regularization. We find that higher task complexity and stronger feature learning reduce degeneracy in neural dynamics but increase it in weight space, with mixed effects on behavior. In contrast, larger networks and structural regularization reduce degeneracy at all three levels. These findings empirically validate the Contravariance Principle and provide practical guidance for researchers aiming to tailor RNN solutions—whether to uncover shared neural mechanisms or to model individual variability observed in biological systems. This work provides a principled framework for quantifying and controlling solution degeneracy in task-trained RNNs, offering new tools for building more interpretable and biologically grounded models of neural computation.

1 Introduction

2

3

5

8

9

10

11

12

13

14

15

16

17

18

19

20

29

30

31

33

34

35

36

Recurrent neural networks (RNNs) are widely used in machine learning and computational neuroscience to model dynamical processes [1, 2, 3, 4, 5, 6]. Traditionally, the study of task-trained RNNs has focused on reverse-engineering a single trained model, implicitly assuming that networks trained on the same task would converge to similar solutions—even when initialized or trained differently. However, recent work has shown that this assumption does not hold universally, and the solution space of task-trained RNNs can be highly degenerate: networks may achieve the same level of training loss, yet differ in out-of-distribution (OOD) behavior, internal representations, neural dynamics, and connectivity [7, 8, 9, 10, 11].

These raise fundamental questions about the solution space of task-trained RNNs: **What factors govern the solution degeneracy across independently trained RNNs?** Despite extensive work in feedforward networks showing how different initializations and stochastic training can yield divergent solutions, RNNs still lack a systematic and unified understanding of the factors that govern solution degeneracy [12, 13, 14, 15, 16, 17, 18, 19, 20]. Cao and Yamins [21] proposed the *Contravariance Principle*, which posits that as the computational objective (i.e., the task) becomes more complex, the solution space should become less dispersed. While this principle is intuitive and compelling, it has thus far remained largely theoretical and has not been directly validated through empirical studies.

In this paper, we introduce a unified framework for quantifying solution degeneracy at three levels: behavior, neural dynamics, and weight space. As illustrated in Figure 1, we quantify degeneracy across behavior, dynamics, and weights, and examine how it is shaped by four key factors. Leveraging this framework, we isolate four key factors that control solution degeneracy—task complexity, learning regime, network width, and structural regularization. By systematically varying task complexity, learning regime, network width, and regularization, we map how each factor shapes degeneracy across behavior, dynamics, and weights.

We find that as task complexity increases—whether via 44 more input-output channels, higher memory demand, or 45 auxiliary objectives—or as networks undergo stronger fea-46 ture learning, their neural dynamics become more consis-47 tent, while their weight configurations grow more variable. 48 In contrast, increasing network size or imposing structural 49 regularization during training reduces variability at both 50 the dynamics and weight levels. At the behavioral level, 51 each of these factors reliably modulates behavioral degen-52 eracy; however, the relationship between behavioral and dynamical degeneracy is not always consistent.



Figure 1: Key factors shape degeneracy across behavior, dynamics, and weights.

2 Methods

55

Model architecture and training. We use discrete-time nonlinear vanilla RNNs with update $\mathbf{h}_t = \tanh\left(\mathbf{W}_h\mathbf{h}_{t-1} + \mathbf{W}_x\mathbf{x}_t + \mathbf{b}\right)$ where $\mathbf{h}_t \in \mathbb{R}^n$ is the hidden state, $\mathbf{x}_t \in \mathbb{R}^m$ is the input, $\mathbf{W}_h \in \mathbb{R}^{n \times n}$ and $\mathbf{W}_x \in \mathbb{R}^{n \times m}$ are the recurrent and input weight matrices. A linear readout maps \mathbf{h}_t to outputs. Networks are trained with BPTT (Adam optimizer, no weight decay) [22]. For each task, we train 50 seeds with 128 hidden units, initializing $\mathbf{W}_h, \mathbf{W}_x \sim \mathcal{U}(-1/\sqrt{n}, 1/\sqrt{n})$. Training continues until networks reach a near-asymptotic training loss threshold, after which we allow 3 epochs' patience period and stop training to assess degeneracy across solutions (Appendix G).

Tasks. We evaluate four neuroscience-relevant tasks eliciting distinct dynamics: pattern recognition (N-Bit Flip-Flop), delayed decision-making (Delayed Discrimination), pattern generation (Sine Wave Generation), and evidence accumulation (Path Integration). Task details and example neural dynamics required to solve the tasks are in Appendix A and F.

Degeneracy metrics. Behavioral degeneracy measures the variability in network responses to out-of-distribution (OOD) inputs. We measure OOD performance as the mean squared error of all converged networks under a temporal generalization condition (double the delay for Delayed Discrimination; double the trial length otherwise). Behavioral degeneracy is the standard deviation of the OOD losses.

Dynamical degeneracy quantifies the average pairwise difference in networks' neural dynamics

71 72 through Dynamical Similarity Analysis (DSA) [23]. DSA compares the topological structure of dynamical systems and has been shown to be more robust to noise and better at identifying behaviorally 73 relevant differences than prior metrics such as Procrustes Analysis and Central Kernel Alignment [24]. 74 For a pair of networks X and Y, DSA identifies a linear forward operator for each system— A_x and 75 A_y —which maps neural activity from one time step to the next. These operators are then compared up 76 to a rotation. The DSA distance between two systems is computed by minimizing the Frobenius norm 77 between the operators, up to rotation: $d_{DSA}(A_x, A_y) = \min_{C \in O(n)} ||A_x - CA_yC^{-1}||_F$, where 78 O(n) is the orthogonal group. We define dynamical degeneracy as the average DSA distance across 79 all network pairs. Additional details are provided in Appendix I. 80

Weight degeneracy is defined as a permutation-invariant Frobenius distance between recurrent weights $d_{\mathrm{PIF}}(\mathbf{W}_1,\mathbf{W}_2) = \min_{\mathbf{P} \in \mathcal{P}(n)} \left\| \mathbf{W}_1 - \mathbf{P}^\top \mathbf{W}_2 \mathbf{P} \right\|_F, \text{ normalized by parameter count when comparing different widths (Appendix I.2).}$

3 Results

84

85

86

87

88

89

3.1 Task complexity modulates degeneracy across levels

We varied task complexity by increasing the number of independent input—output channels of each task, which effectively duplicated the task across dimensions and increased the representational load of networks by forcing them to multitask. Higher task complexity constrains the space of viable dynamical solutions, leading to tighter clustering and greater similarity across independently trained networks (Fig. 2AB). At the behavioral level, networks trained on more complex tasks consistently showed greater consistency and lower variability in their responses to OOD test inputs

(Fig 2D). Together, the results at the behavioral and dynamical levels support the *Contravariance Principle*, which posits an inverse relationship between task complexity and the dispersion of network solutions [21]. At the weight level, however, we found that pairwise distances between converged RNNs' weight matrices increased consistently with task complexity (Figure 2C), which likely reflects increased dispersion of local minima in weight space for harder tasks [25, 26, 27, 28, 29, 30, 31]. In Appendix B, we explore two alternative approaches of varying task complexity: increasing the task's memory demand and adding auxiliary objectives. We find that the trends in solution degeneracy hold consistently across these approaches.

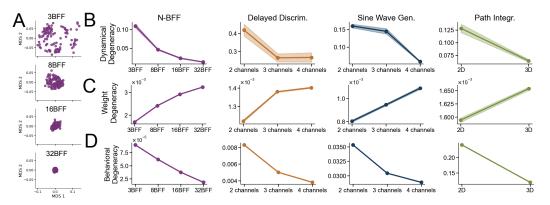


Figure 2: Higher task complexity reduces dynamical and behavioral degeneracy, but increases weight degeneracy. (A) Two-dimensional MDS embedding of network dynamics shows that independently trained networks converge to more similar trajectories as task complexity increases. (B) Dynamical, (C) weight, and (D) behavioral degeneracy across 50 networks as a function of task complexity. Shaded area indicates ± 1 standard error.

3.2 Controlling feature learning reshapes degeneracy across levels

In deep learning theory, neural networks can operate in either a lazy or rich learning regime [32, 33, 34, 35]. In the lazy regime, weights and internal features remain largely unchanged during training. In the rich (feature learning) regime, networks reshape their hidden representations and weights to capture task-specific structure [32, 36, 37, 33].

Intuitively, when networks undergoes strong feature learning, they converge to more consistent task-specific neural dynamics, leading to lower dynamical degeneracy. causally test whether feature learning affects solution degeneracy, we used a principled parameterization known as maximum update parameterization (μP) , where a single hyperparameter— γ —controls the strength of feature learning: higher γ values induce richer feature-learning regime [35, 32, 34, 33]. More specifically, the network output is scaled as

93

94

95

96

97

100

101

102

103

104

105

106

107

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

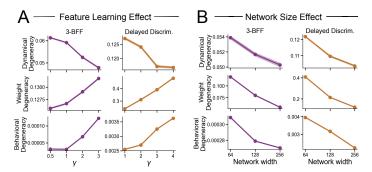


Figure 3: (A) Stronger feature learning reduces dynamical degeneracy but increases weight and behavioral degeneracy. (B) Larger networks reduce degeneracy across weight, dynamics, and behavior. Panels show degeneracy at the dynamical, weight, and behavioral levels (top to bottom). Shaded area indicates ± 1 standard error.

 $f(t) = \frac{1}{\gamma N} W_{\text{readout}} \phi(h(t))$. A detailed explanation of μP and its relationship to the standard parameterization is in Appendix L and M. For each task, we trained networks with multiple γ values and confirmed that larger γ consistently induces stronger feature learning(Appendix N).

We observed that stronger feature learning reduced degeneracy at the dynamical level but increased it at the weight level (consistent for all four tasks, see Appendix D). This finding aligns with prior

work in feedforward networks, where feature learning was shown to reduce the variance of the neural tangent kernel across converged models [38]. Notably, stronger feature learning was shown to push networks to travel farther from their initialization [39, 36], resulting in more dispersed final weights and higher weight degeneracy. At the behavioral level, however, increasing feature-learning strength leads networks to overfit the training distribution (Appendix K.2). We hypothesize that stronger feature learning exacerbates overfitting, increasing both average OOD loss and the variability of OOD behavior across models (Figure 8) [40, 41, 42, 43].

3.3 Larger networks yield more consistent solutions across levels

Although larger networks may yield more consistent solutions via self-averaging and improved convergence [44, 45, 46, 47, 48], this outcome is not guaranteed without controlling for feature learning, as increasing network width pushes models towards the lazy regime, where feature learning is suppressed [49, 37, 32, 33, 34]. To disentangle these competing effects, we again use the μP parameterization, which holds feature learning strength constant (via fixed γ) while scaling width. Across all tasks, larger networks consistently exhibit lower degeneracy at the weight, dynamical, and behavioral levels, producing more consistent solutions across random seeds (Figure 9; results hold consistent for all four tasks, see Appendix E). This pattern aligns with findings in vision and language models, where wider networks converge to more similar internal representations [50, 51, 52, 53, 54].

3.4 Structural regularization reduces solution degeneracy

Low-rank and sparsity constraints are widely used structural regularizers in neuroscience-inspired modeling and efficient machine learning [4, 55, 56, 57, 58]. A low-rank penalty compresses the weight matrices into a few dominant modes, while an ℓ_1 penalty drives many parameters to zero and induces sparsity. In both cases, task-irrelevant features are pruned, nudging independently initialized networks toward more consistent solutions on the same task. To test this idea, we augmented the task loss with either a nuclear-norm penalty on the recurrent weights $\mathcal{L} = \mathcal{L}_{\text{task}} + \lambda_{\text{rank}} \sum_{i=1}^{r} \sigma_i$, where σ_i are the singular values of the recurrent matrix, or an ℓ_1 sparsity penalty: $\mathcal{L} = \mathcal{L}_{\text{task}} + \lambda_{\ell_1} \sum_i |w_i|$. We focused on the Delayed Discrimination task to control for baseline difficulty, but both regularizers consistently reduced degeneracy across all levels—and similar effects hold in other tasks (Appendix P, Figure 4).

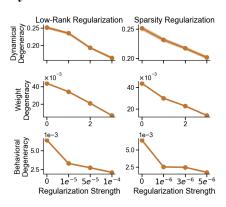


Figure 4: Low-rank and sparsity regularization reduce solution degeneracy across all levels. Shaded area indicates ± 1 standard error.

4 Discussion

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

166

167

168

169

Table 1: Summary of how each factor affects solution degeneracy. Arrows indicate the direction of change for each level as the factor increases. Contravariant factors shift **dynamic** and **weight degeneracy** in opposite direction; covariant factors shift them in the same directions.

Factor	Dynamics	Weights	Behavior
Higher Task complexity (contravariant)	+	<u></u>	
More Feature learning (contravariant)	↓	1	1
Larger Network size (covariant) Regularization (covariant)	‡	‡	‡

We present a unified framework for quantifying solution degeneracy in task-trained RNNs, identify the key factors that shape the solution landscape. In both machine learning and neuroscience, the optimal level of degeneracy may vary depending on the specific research questions being investigated. This framework offers practical guidance for tailoring training to a given goal—whether encouraging consistency across models [59], or promoting diversity across learned solutions [60, 61, 62].

References

170

- [1] David Sussillo. Neural circuits as computational dynamical systems. *Current opinion in neurobiology*, 25:156–163, 2014.
- [2] Kanaka Rajan, Christopher D Harvey, and David W Tank. Recurrent network models of sequence generation and memory. *Neuron*, 90(1):128–142, 2016.
- [3] Omri Barak. Recurrent neural networks as versatile tools of neuroscience research. 46:1–6. ISSN 09594388. doi: 10.1016/j.conb.2017.06.003. URL https://linkinghub.elsevier.com/retrieve/pii/S0959438817300429.
- 178 [4] Francesca Mastrogiuseppe and Srdjan Ostojic. Linking connectivity, dynamics, and computations in low-rank recurrent neural networks. *Neuron*, 99(3):609–623.e29, 2018. doi: 10.1016/j.neuron.2018.07.003.
- [5] Saurabh Vyas, Matthew D Golub, David Sussillo, and Krishna V Shenoy. Computation through
 neural population dynamics. *Annual Review of Neuroscience*, 43:249–275, 2020.
- [6] Laura N. Driscoll, Krishna Shenoy, and David Sussillo. Flexible multitask computation in recurrent networks utilizes shared dynamical motifs. *Nature Neuroscience*, 27(7):1349–1363,
 July 2024. ISSN 1097-6256, 1546-1726. doi: 10.1038/s41593-024-01668-6. URL https://www.nature.com/articles/s41593-024-01668-6.
- [7] Elia Turner, Kabir V Dabholkar, and Omri Barak. Charting and navigating the space of
 solutions for recurrent neural networks. Advances in Neural Information Processing Systems,
 34:25320–25333, 2021.
- [8] Aniruddh R Galgali, Maneesh Sahani, and Valerio Mante. Residual dynamics resolves recurrent contributions to neural computation. *Nature Neuroscience*, 26(2):326–338, 2023.
- [9] Gita Gholamrezaei and Ian Q. Whishaw. Individual differences in skilled reaching for food related to increased number of gestures: Evidence for goal and habit learning of skilled reaching. 123(4):863–874. ISSN 1939-0084, 0735-7044. doi: 10.1037/a0016369. URL http://doi.apa.org/getdoi.cfm?doi=10.1037/a0016369.
- [10] Peiran Gao, Eric Trautmann, Byron Yu, Gopal Santhanam, Stephen Ryu, Krishna Shenoy, and
 Surya Ganguli. A theory of multineuronal dimensionality, dynamics and measurement. URL
 http://biorxiv.org/lookup/doi/10.1101/214262.
- [11] Johannes Mehrer, Courtney J. Spoerer, Nikolaus Kriegeskorte, and Tim C. Kietzmann. Individual differences among deep neural network models. 11(1):5725. ISSN 2041-1723. doi: 10.1038/s41467-020-19632-w. URL http://www.nature.com/articles/s41467-020-19632-w.
- 202 [12] Abhranil Das and Ila R Fiete. Systematic errors in connectivity inferred from activity in strongly recurrent networks. *Nature Neuroscience*, 23(10):1286–1296, 2020.
- 204 [13] Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, pages 249–256. JMLR Workshop and Conference Proceedings, 2010.
- 207 [14] Stanislav Fort, Huiyi Hu, and Balaji Lakshminarayanan. Deep ensembles: A loss landscape perspective. *arXiv preprint arXiv:1912.02757*, 2019.
- 209 [15] Ian J. Goodfellow, Oriol Vinyals, and Andrew M. Saxe. Qualitatively characterizing neural 210 network optimization problems, 2015. URL https://arxiv.org/abs/1412.6544.
- [16] Hao Li, Zheng Xu, Gavin Taylor, Christoph Studer, and Tom Goldstein. Visualizing the loss
 landscape of neural nets, 2018. URL https://arxiv.org/abs/1712.09913.
- 213 [17] Stanisław Jastrzębski, Zachary Kenton, Devansh Arpit, Nicolas Ballas, Asja Fischer, Yoshua 214 Bengio, and Amos Storkey. Three factors influencing minima in sgd, 2018. URL https: 215 //arxiv.org/abs/1711.04623.

- 216 [18] Pratik Chaudhari, Anna Choromanska, Stefano Soatto, Yann LeCun, Carlo Baldassi, Christian Borgs, Jennifer Chayes, Levent Sagun, and Riccardo Zecchina. Entropy-sgd: Biasing gradient descent into wide valleys, 2017. URL https://arxiv.org/abs/1611.01838.
- [19] Jonathan Frankle and Michael Carbin. The lottery ticket hypothesis: Finding sparse, trainable neural networks, 2019. URL https://arxiv.org/abs/1803.03635.
- 221 [20] Simon Kornblith, Mohammad Norouzi, Honglak Lee, and Geoffrey Hinton. Similarity of neural network representations revisited, 2019. URL https://arxiv.org/abs/1905.00414.
- 223 [21] Rosa Cao and Daniel Yamins. Explanatory models in neuroscience, part 2: Functional intelligibility and the contravariance principle. *Cognitive Systems Research*, 85:101200, 2024.
- ²²⁵ [22] Paul J Werbos. Backpropagation through time: what it does and how to do it. *Proceedings of the IEEE*, 78(10):1550–1560, 1990.
- 227 [23] Mitchell Ostrow, Adam Eisen, Leo Kozachkov, and Ila Fiete. Beyond Geometry: Comparing 228 the Temporal Structure of Computation in Neural Circuits with Dynamical Similarity Analysis, 229 October 2023. URL http://arxiv.org/abs/2306.10168. arXiv:2306.10168 [cs, q-bio].
- 230 [24] Quentin Guilhot, Michał J Wójcik, Jascha Achterberg, and Rui Ponte Costa. Dynamical similarity analysis uniquely captures how computations develop in RNNs, 2025. URL https://openreview.net/forum?id=pXPIQsV1St.
- [25] Ian J. Goodfellow, Oriol Vinyals, and Andrew M. Saxe. Qualitatively characterizing neural
 network optimization problems. In *International Conference on Learning Representations* (ICLR), 2015. arXiv:1412.6544.
- [26] Jonathan Frankle, Gintare Karolina Dziugaite, Daniel Roy, and Michael Carbin. Linear mode
 connectivity and the lottery ticket hypothesis. In *Proceedings of the 37th International Conference on Machine Learning (ICML)*, pages 3259–3269. PMLR, 2020.
- [27] James R. Lucas, Juhan Bae, Michael R. Zhang, Stanislav Fort, Richard Zemel, and Roger B.
 Grosse. On monotonic linear interpolation of neural network parameters. In *Proceedings of the* 38th International Conference on Machine Learning (ICML), pages 7168–7179. PMLR, 2021.
- 242 [28] Stanislav Fort and Stanislaw Jastrzebski. Large scale structure of neural network loss landscapes.
 243 In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 32, 2019.
- 244 [29] Alessandro Achille, Giovanni Paolini, and Stefano Soatto. Where is the information in a deep neural network? *CoRR*, abs/1905.12213, 2019.
- [30] Xingyu Qu and Samuel Horvath. Rethink model re-basin and the linear mode connectivity.
 arXiv preprint arXiv:2402.05966, 2024.
- 248 [31] Andrew Ly and Pulin Gong. Optimization on multifractal loss landscapes explains a diverse 249 range of geometrical and dynamical properties of deep learning. *Nature Communications*, 16 250 (3252), 2025. doi: 10.1038/s41467-025-58532-9.
- [32] Léon Chizat, Edouard Oyallon, and Francis Bach. On lazy training in differentiable programming.
 Advances in Neural Information Processing Systems, 32:2938–2950, 2019.
- [33] Bryan Woodworth, Suriya Gunasekar, Jason D. Lee, Nathan Srebro, Srinadh Bhojanapalli, Rina
 Khanna, Aaron Chatterji, and Martin Jaggi. Kernel and rich regimes in deep learning. *Journal* of Machine Learning Research, 21(243):1–48, 2020.
- [34] Mario Geiger, Stefano Spigler, Arthur Jacot, and Matthieu Wyart. Disentangling feature and
 lazy training in deep neural networks. *Journal of Statistical Mechanics: Theory and Experiment*,
 2020(11):113301, 2020. doi: 10.1088/1742-5468/abc4de.
- 259 [35] Blake Bordelon and Cengiz Pehlevan. Self-Consistent Dynamical Field Theory of Kernel 260 Evolution in Wide Neural Networks, October 2022. URL http://arxiv.org/abs/2205. 261 09653. arXiv:2205.09653 [stat].

- ²⁶² [36] Thomas George, Guillaume Lajoie, and Aristide Baratin. Lazy vs hasty: Linearization in deep networks impacts learning schedule based on example difficulty. *arXiv* preprint *arXiv*:2209.09658, 2022. URL https://arxiv.org/abs/2209.09658.
- Jaehoon Lee, Yuval Bahri, Roman Novak, Samuel S. Schoenholz, Jeffrey Pennington, and
 Jascha Sohl-Dickstein. Wide neural networks of any depth evolve as linear models under
 gradient descent. In Advances in Neural Information Processing Systems, volume 32, pages
 8572–8583, 2019.
- 269 [38] Blake Bordelon and Cengiz Pehlevan. Dynamics of finite width Kernel and prediction fluctua-270 tions in mean field neural networks*. *Journal of Statistical Mechanics: Theory and Experiment*, 271 2024(10):104021, October 2024. ISSN 1742-5468. doi: 10.1088/1742-5468/ad642b. URL 272 https://iopscience.iop.org/article/10.1088/1742-5468/ad642b.
- Yuhan Helena Liu, Aristide Baratin, Jonathan Cornford, Stefan Mihalas, Eric Shea-Brown, and Guillaume Lajoie. How connectivity structure shapes rich and lazy learning in neural circuits. arXiv preprint arXiv:2310.08513, 2023. doi: 10.48550/arXiv.2310.08513. URL https://arxiv.org/abs/2310.08513.
- 277 [40] Yamini Bansal, Preetum Nakkiran, and Boaz Barak. Revisiting model stitching to compare neural representations. *arXiv preprint arXiv:2106.07682*, 2021. URL https://arxiv.org/abs/2106.07682.
- 280 [41] Sunny Duan, Loïc Matthey, André Saraiva, Nicholas Watters, Christopher P. Burgess, Alexander Lerchner, and Irina Higgins. Unsupervised model selection for variational disentangled representation learning. *arXiv preprint arXiv:1905.12614*, 2020. URL https://arxiv.org/abs/1905.12614.
- [42] Minyoung Huh, Brian Cheung, Tongzhou Wang, and Phillip Isola. The platonic representation hypothesis. arXiv preprint arXiv:2405.07987, 2024. URL https://arxiv.org/abs/2405.07987.
- 287 [43] Yixuan Li, Jason Yosinski, Jeff Clune, Hod Lipson, and John Hopcroft. Convergent learning:
 288 Do different neural networks learn the same representations? *arXiv preprint arXiv:1511.07543*,
 289 2016. URL https://arxiv.org/abs/1511.07543.
- 290 [44] Kenji Kawaguchi. Deep learning without poor local minima. In *Advances in Neural Information* 291 *Processing Systems* 29, pages 586–594, 2016.
- Quynh Nguyen and Matthias Hein. The loss surface of deep and wide neural networks. *CoRR*, abs/1704.08045, 2017.
- 294 [46] Simon S. Du, Jason D. Lee, Haochuan Li, Liwei Wang, and Xiyu Zhai. Gradient descent finds 295 global minima of deep neural networks. *CoRR*, abs/1811.03804, 2018.
- Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song. A convergence theory for deep learning
 via over-parameterization. In *Proceedings of the 36th International Conference on Machine Learning*, pages 242–252, 2019.
- [48] Difan Zou, Yuan Cao, Dongruo Zhou, and Quanquan Gu. Stochastic gradient descent optimizes
 over-parameterized deep relu networks. *CoRR*, abs/1811.08888, 2018.
- [49] Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization in neural networks. In *Advances in Neural Information Processing Systems*, volume 31, 2018.
- [50] Ari S. Morcos, Maithra Raghu, and Samy Bengio. Insights on representational similarity inneural networks with canonical correlation. In *NeurIPS*, 2018.
- [51] Simon Kornblith, Mohammad Norouzi, Honglak Lee, and Geoffrey Hinton. Similarity of neural
 network representations revisited. In *ICML*, 2019.
- Maithra et al. Raghu. Svcca: Singular vector canonical correlation analysis for deep learning dynamics. In *NeurIPS*, 2017.

- [53] Fred Wolf, Rainer Engelken, Maximilian Puelma-Touzel, Juan Daniel Flórez Weidinger, and
 Andreas Neef. Dynamical models of cortical circuits. 25:228–236. ISSN 09594388. doi:
 10.1016/j.conb.2014.01.017. URL https://linkinghub.elsevier.com/retrieve/pii/
 S0959438814000324.
- Felix et al. Klabunde. Contrasim analyzing neural representations based on contrastive learning. In *ICLR*, 2024.
- [55] Manuel Beiran, Alexis Dubreuil, Adrian Valente, Francesca Mastrogiuseppe, and Srdjan Ostojic.
 Shaping dynamics with multiple populations in low-rank recurrent networks. *arXiv preprint arXiv:2007.02062*, 2020. doi: 10.48550/arXiv.2007.02062.
- [56] Bruno A. Olshausen and David J. Field. Emergence of simple-cell receptive field properties by
 learning a sparse code for natural images. *Nature*, 381(6583):607–609, 1996. doi: 10.1038/321
 381607a0.
- Song Han, Jeff Pool, John Tran, and William J. Dally. Learning both weights and connections for efficient neural networks. In *Advances in Neural Information Processing Systems*, pages 1135–1143, 2015. doi: 10.48550/arXiv.1506.02626.
- [58] Xavier Glorot, Antoine Bordes, and Yoshua Bengio. Deep sparse rectifier neural networks.
 In Proceedings of the 14th International Conference on Artificial Intelligence and Statistics,
 volume 15, pages 315–323, 2011.
- 528 [59] D Kepple, Rainer Engelken, and Kanaka Rajan. Curriculum learning as a tool to uncover learning principles in the brain. In *International Conference on Learning Representations*, 2022.
- [60] Samuel Liebana Garcia, Aeron Laffere, Chiara Toschi, Louisa Schilling, Jacek Podlaski,
 Matthias Fritsche, Peter Zatka-Haas, Yulong Li, Rafal Bogacz, Andrew Saxe, and Armin
 Lak. Striatal dopamine reflects individual long-term learning trajectories, December 2023. URL
 http://biorxiv.org/lookup/doi/10.1101/2023.12.14.571653.
- Valeria Fascianelli, Aldo Battista, Fabio Stefanini, Satoshi Tsujimoto, Aldo Genovesio, and
 Stefano Fusi. Neural representational geometries reflect behavioral differences in monkeys and
 recurrent neural networks. *Nature Communications*, 15(1):6479, August 2024. ISSN 2041 doi: 10.1038/s41467-024-50503-w. URL https://www.nature.com/articles/s41467-024-50503-w.
- [62] A Pan-Vazquez, Y Sanchez Araujo, B McMannon, M Louka, A Bandi, L Haetzel, International
 Brain Laboratory, JW Pillow, ND Daw, and IB Witten. Pre-existing visual responses in a
 projection-defined dopamine population explain individual learning trajectories, 2024. URL
 https://europepmc.org/article/PPR/PPR811803.
- ³⁴³ [63] Peter J Schmid. Dynamic mode decomposition and its variants. *Annual Review of Fluid Mechanics*, 54(1):225–254, 2022.
- ³⁴⁵ [64] Peter H. Schönemann. A generalized solution of the orthogonal procrustes problem. *Psychometrika*, 31(1):1–10, Mar 1966. doi: 10.1007/BF02289451.
- [65] Chris Ding, Tao Li, and Michael I. Jordan. Nonnegative matrix factorization for combinatorial optimization: Spectral clustering, graph matching, and clique finding. In *Proceedings of the Eighth IEEE International Conference on Data Mining (ICDM '08)*, pages 183–192. IEEE, 2008. doi: 10.1109/ICDM.2008.130.
- [66] Fanwang Meng, Michael G. Richer, Alireza Tehrani, Jonathan La, Taewon David Kim, P. W. Ayers, and Farnaz Heidar-Zadeh. Procrustes: A python library to find transformations that maximize the similarity between matrices. *Computer Physics Communications*, 276:108334, 2022. doi: 10.1016/j.cpc.2022.108334. URL https://www.sciencedirect.com/science/article/pii/S0010465522000522.
- Greg Yang, Edward J. Hu, Igor Babuschkin, Szymon Sidor, Xiaodong Liu, David Farhi, Nick Ryder, Jakub Pachocki, Weizhu Chen, and Jianfeng Gao. Tensor programs v: Tuning large neural networks via zero-shot hyperparameter transfer. *arXiv preprint arXiv:2203.03466*, 2022. doi: 10.48550/arXiv.2203.03466. Accepted at NeurIPS 2021.

- [68] Greg Yang, Dingli Yu, Chen Zhu, and Soufiane Hayou. Tensor programs vi: Feature learning in infinite-depth neural networks. *arXiv preprint arXiv:2310.02244*, 2023. doi: 10.48550/arXiv. 2310.02244. Accepted at ICLR 2024.
- Kenneth D Miller and Francesco Fumarola. Mathematical equivalence of two common forms of firing rate models of neural networks. *Neural computation*, 24(1):25–31, 2012.

365 A Task descriptions

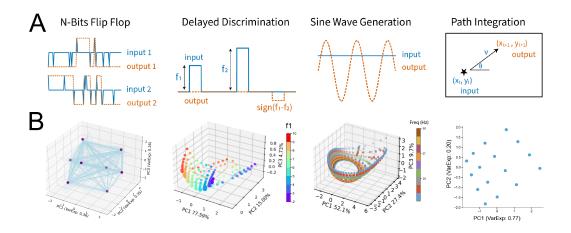


Figure 5: Our task suite spans memory, integration, pattern generation, and decision-making. Each task is designed to place distinct demands on the network's dynamics. N-Bit Flip-Flop: The network must remember the last nonzero input on each of N independent channels. Delayed Discrimination: The network compares the magnitude of two pulses, separated by a variable delay, and outputs their sign difference. Sine Wave Generation: A static input specifies a target frequency, and the network generates the corresponding sine wave over time. Path Integration: The network integrates velocity inputs to track position in a bounded 2D or 3D arena (schematic shows 2D case).

N-Bit Flip-Flop Task Each RNN receives N independent input channels taking values in $\{-1,0,+1\}$, which switch with probability p_{switch} . The network has N output channels that must retain the most recent nonzero input on their respective channels. The network dynamics form 2^N fixed points, corresponding to all binary combinations of $\{-1,+1\}^N$.

Delayed Discrimination Task The network receives two pulses of amplitudes $f_1, f_2 \in [2, 10]$, separated by a variable delay $t \in [5, 20]$ time steps, and must output $\operatorname{sign}(f_2 - f_1)$. In the N-channel variant, comparisons are made independently across channels. The network forms task-relevant fixed points to retain the amplitude of f_1 during the delay period.

Sine Wave Generation The network receives a static input specifying a target frequency $f \in [1,30]$ and must generate the corresponding sine wave $\sin(2\pi ft)$ over time. We define N_{freq} target frequencies, evenly spaced within the range [1,30], and use them during training. In the N-channel variant, each input channel specifies a frequency, and the corresponding output channel generates a sine wave at that frequency. For each frequency, the network dynamics form and traverse a limit cycle that produces the corresponding sine wave.

Path Integration Task Starting from a random position in 2D, the network receives angular direction θ and speed v at each time step and updates its position estimate. In the 3D variant, the network takes as input azimuth θ , elevation ϕ , and speed v, and outputs updated (x,y,z) position. The network performs path integration by accumulating velocity vectors based on the input directions and speeds. After training, the network forms a Euclidean map of the environment in its internal state space.

B Additional axes of task complexity

In the main text, we controlled task complexity by varying the number of independent input—output channels, effectively duplicating the task across dimensions. Here, we explore two alternative approaches: increasing the task's memory demand and adding auxiliary objectives.

Changing memory demand. Of the four tasks, only Delayed Discrimination requires extended memory, as its performance depends on maintaining the first stimulus across a variable delay. See Appendix H for a quantification of each task's memory demand. We increased the memory load in Delayed Discrimination by lengthening the delay period. This manipulation reduced degeneracy

at the dynamical and behavioral levels but increased it at the weight level, mirroring the effect of increasing task dimensionality (Figure 6A).

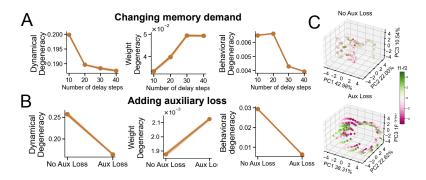


Figure 6: Memory demand and auxiliary loss modulate degeneracy in distinct ways. In the Delayed Discrimination task, both manipulations reduce dynamical and behavioral degeneracy while increasing weight degeneracy. The auxiliary loss also induces additional line attractors in the network's dynamics, as shown in (C).

Adding auxiliary loss. We next examined how adding an auxiliary loss affects solution degeneracy in the Delayed Discrimination task. Specifically, the network outputs both the sign and the magnitude of the difference between two stimulus values (f_2-f_1) , using separate output channels for each. This manipulation added a second output channel and increased memory demand by requiring the network to track the magnitude of the difference between incoming stimuli. Consistent with our hypothesis, this manipulation reduced dynamical and behavioral degeneracy while increasing weight degeneracy (Figure 6B). Crucially, the auxiliary loss induced additional line attractors in the network dynamics, further structuring internal trajectories and aligning neural responses across networks (Figure 6C). While the auxiliary loss increases both output dimensionality and temporal memory demand, we interpret its effect holistically as a structured increase in task complexity.

C Higher task complexity induces more feature learning

395

396

397

398

399

400

401

402

403

404

405

We hypothesize that the increased weight degeneracy observed in harder tasks reflects stronger feature 406 learning. Specifically, harder tasks may force network weights to travel farther from their initialization. 407 If more complex task variants, like those in Section 3.1, truly induce greater feature learning, then 408 networks should traverse a greater distance in weight space, resulting in more dispersed final weights. 409 To test this idea, we measured feature learning strength in networks trained on different task variants 410 using two complementary metrics [39, 36]: Weight-change norm: $\|\mathbf{W}_T - \mathbf{W}_0\|_F$, where larger 411 values indicate stronger feature learning. Kernel alignment (KA): measures the directional change of 412 the neural tangent kernel (NTK) before and after training: $KA(K^{(f)}, K^{(0)}) = \frac{Tr(K^{(f)}K^{(0)})}{\|K^{(f)}\|_F \|K^{(0)}\|_F}$ 413 where $K = \nabla_W \hat{y}^\top \nabla_W \hat{y}$. Lower KA indicates greater NTK rotation and thus stronger feature learning. 415 More complex tasks consistently drive stronger feature learning and greater dispersion in weight 416 space, as reflected by increasing weight-change norm and decreasing kernel alignment across all 417 tasks (Figure 7). 418

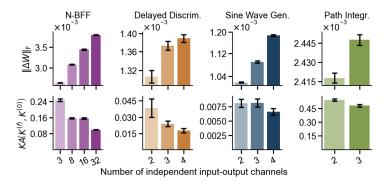


Figure 7: More complex tasks drive stronger feature learning in RNNs. Increased input—output dimensionality leads to higher weight-change norms and lower kernel alignment. Error bars indicate ± 1 standard error.

19 D Feature learning effect for all tasks

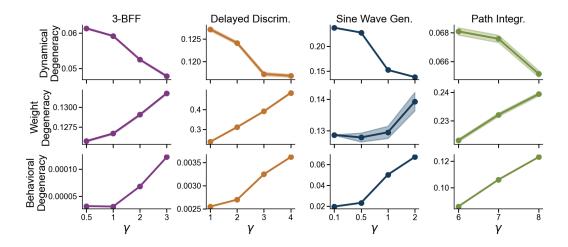


Figure 8: Stronger feature learning reduces dynamical degeneracy but increases weight and behavioral degeneracy. Panels show degeneracy at the dynamical, weight, and behavioral levels (top to bottom). Shaded area indicates ± 1 standard error.

E Network size effect for all tasks

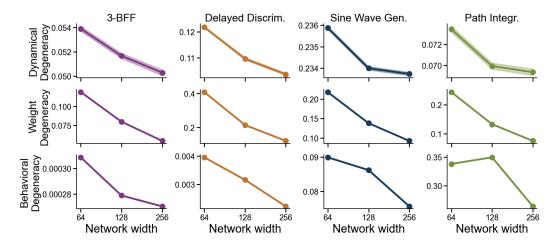


Figure 9: Larger networks reduce degeneracy across weight, dynamics, and behavior. Controlling for feature learning strength, wider RNNs yield more consistent solutions across all three levels of analysis. Panels show degeneracy at the dynamical, weight, and behavioral levels (top to bottom). Shaded area indicates ± 1 standard error.

421 F Task details

422 F.1 N-Bit Flip Flop

Task Parameter	Value
Probability of flip	0.3
Number of time steps	100

423 F.2 Delayed Discrimination

Task Parameter	Value
Number of time steps	60
Max delay	20
Lowest stimulus value	2
Highest stimulus value	10

424 F.3 Sine Wave Generation

Task Parameter	Value
Number of time steps	100
Time step size	0.01
Lowest frequency	1
Highest frequency	30
Number of frequencies	100

425 F.4 Path Integration

Task Parameter	Value
Number of time steps	100
Maximum speed (v_{max})	0.4
Direction increment std ($\theta_{\rm std}$ / $\phi_{\rm std}$)	$\pi/10$
Speed increment std	0.1
Noise std	0.0001
Mean stop duration	30
Mean go duration	50
Environment size (per side)	10

426 G Training details

427 G.1 N-Bit Flip Flop

Training Hyperparameter	Value
Optimizer	Adam
Learning rate	0.001
Learning rate scheduler	None
Max epochs	300
Steps per epoch	128
Batch size	256
Early stopping threshold	0.001
Patience	3
Time constant (μP)	1

428 G.2 Delayed Discrimination

Training Hyperparameter	Value
Optimizer	Adam
Learning rate	0.001
Learning rate scheduler	CosineAnnealingWarmRestarts
Max epochs	500
Steps per epoch	128
Batch size	256
Early stopping threshold	0.01
Patience	3
Time constant (μP)	0.1

429 G.3 Sine Wave Generation

Training Hyperparameter	Value
Optimizer	Adam
Learning rate	0.0005
Learning rate scheduler	None
Max epochs	500
Steps per epoch	128
Batch size	32
Early stopping threshold	0.05
Patience	3
Time constant (μP)	1

430 G.4 Path Integration

Training Hyperparameter	Value
Optimizer	Adam
Learning rate	0.001
Learning rate scheduler	ReduceLROnPlateau
Learning rate decay factor	0.5
Learning rate decay patience	40
Max epochs	1000
Steps per epoch	128
Batch size	64
Early stopping threshold	0.05
Patience	3
Time constant (μP)	0.1

431 H Memory demand of each task

In this section, we quantify each task's memory demand by measuring how far back in time its inputs influence the next output. Specifically, for each candidate history length h, we build feature vectors

$$\mathbf{s}_{t}^{(h)} = [x_{t-h+1}, \dots, x_{t}; y_{t}] \in \mathbb{R}^{h d_{\text{in}} + d_{\text{out}}},$$

and train a two-layer MLP to predict the subsequent target y_{t+1} . We then evaluate the held-out mean-squared error $\mathrm{MSE}(h)$, averaged over multiple random initializations. We identify the smallest history length h^* at which the error curve plateaus or has a minimum, and take h^* as the task's

437 intrinsic memory demand.

From the results, we can see that the N-Bits Flip-Flop task requires only one time-step of mem-438 ory—exactly what's needed to recall the most recent nonzero input in each channel. The Sine Wave 439 Generation task demands two time-steps, reflecting the need to track both phase and direction of 440 change. Path Integration likewise only needs one time-step, since the current position plus instanta-441 neous velocity and heading suffice to predict the next position. Delayed Discrimination is the only 442 memory-intensive task: our method estimates a memory demand of 25 time-steps, which happens 443 to be the time interval between the offset of the first stimulus and the onset of the response period, during which the network needs to first keep track of the amplitude of the first stimulus and then its 445 decision. 446

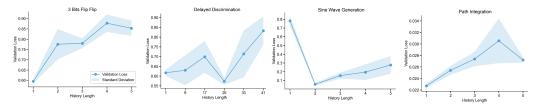


Figure 10: **Memory demand of each task**. The held-out mean-squared error MSE(h) of a two-layer MLP predictor is plotted against history length h. The intrinsic memory demand h^* , defined by the plateau or minimum of each curve, is 1 for the N-Bits Flip-Flop and Path Integration tasks, 2 for Sine Wave Generation, and 25 for Delayed Discrimination—matching the inter-stimulus delay interval in that task.

47 I More details on the degeneracy metrics

448 I.1 Dynamical Degeneracy

Briefly, DSA proceeds as follows: Given two RNNs with hidden states $\mathbf{h}_1(t) \in \mathbb{R}^n$ and $\mathbf{h}_2(t) \in \mathbb{R}^n$, we first generate a delay-embedded matrix, \mathbf{H}_1 and \mathbf{H}_2 of the hidden states in their original state space. Next, for each delay-embedded matrix, we use Dynamic Mode Decomposition (DMD) [63] to extract linear forward operators \mathbf{A}_1 and \mathbf{A}_2 of the two systems' dynamics. Finally, a Procrustes distance between the two matrices \mathbf{A}_1 and \mathbf{A}_2 is used to quantify the dissimilarity between the two dynamical systems and provide an overall DSA score, defined as:

$$d_{\text{Procrustes}}(\mathbf{A}_1, \mathbf{A}_2) = \min_{\mathbf{Q} \in O(n)} \|\mathbf{A}_1 - \mathbf{Q}\mathbf{A}_2\mathbf{Q}^{-1}\|_F$$

where \mathbf{Q} is a rotation matrix from the orthogonal group O(n) and $\|\cdot\|_F$ is the Frobenius norm. This metric quantifies how dissimilar the dynamics of the two RNNs are after accounting for orthogonal transformations. We quantify Dynamical Degeneracy across many RNNs as the average pairwise distance between pairs of RNN neural-dynamics (hidden-state trajectories).

After training, we extract each network's hidden-state activations for every trial in the training set, 459 yielding a tensor of shape (trials \times time steps \times neurons). We collapse the first two dimensions and 460 yield a matrix of size (trials \times time steps) \times neurons. We then apply PCA to retain the components 461 that explain 99% of the variance to remove noisy and low-variance dimensions of the hidden state 462 trajectories. Next, we perform a grid search over candidate delay lags, with a minimum lag of 1 463 and a maximum lag of 30, selecting the lag that minimizes the reconstruction error of DSA on the 464 dimensionality reduced trajectories. Finally, we fit DSA with full rank and the optimal lag to these 465 466 PCA-projected trajectories and compute the pairwise DSA distances between all networks.

I.2 Weight degeneracy

467

468

469

We computed the pairwise distance between the recurrent matrices from different networks using Two-sided Permutation with One Transformation [64, 65] function from the Procrustes Python package [66].

471 J Representational degeneracy

We further quantified solution degeneracy at the representational level—that is, the variability in each network's internal feature space when presented with the same input dataset—using Singular Vector Canonical Correlation Analysis (SVCCA). SVCCA works by first applying singular value decomposition (SVD) to each network's activation matrix, isolating the principal components that capture most of its variance, and then performing canonical correlation analysis (CCA) to find the maximally correlated directions between the two reduced subspaces. The resulting canonical correlations therefore measure how similarly two networks represent the same inputs: high average correlations imply low representational degeneracy (i.e., shared feature subspaces), whereas lower correlations reveal greater divergence in what the models learn. We define the representational degeneracy (labeled as the SVCCA distance below) as

$$d_{\text{repr}}(A_x, A_y) = 1 - \text{SVCCA}(A_x, A_y).$$

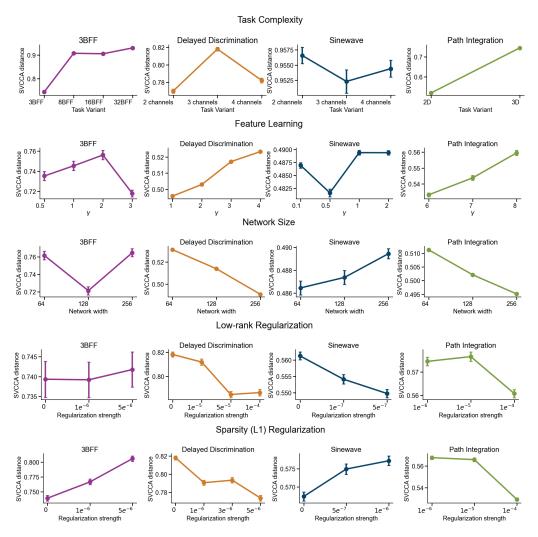


Figure 11: Representational degeneracy, as measured by the average SVCCA distance between networks, does not necessarily change uniformly as we vary task complexity, feature learning strength, network size, and regularization strength.

We found that as we vary the four factors that robustly control the dynamical degeneracy across task-trained RNNs, the representational-level degeneracy isn't necessarily constrained by those same factors in the same way. In RNNs, task-relevant computations are implemented at the level of

network's dynamics instead of static representations, and RNNs that implement similar temporal dynamics can have disparate representational geometry. Therefore, it is expected that task complexity, learning regime, and network size change the task-relevant computations learned by the networks by affecting their neural dynamics instead of representations. DSA captures the dynamical aspect of the neural computation by fitting a forward operator matrix *A* that maps the network's activity at one time step to the next, therefore directly capturing the temporal evolution of neural activities. By contrast, SVCCA aligns the principal subspaces of activation vectors at each time point but treats those vectors as independent samples—it never examines how one state evolves into the next. As a result, SVCCA measures only static representational similarity and cannot account for the temporal dependencies that underlie RNN computations. Nonetheless, we expect SVCCA might be more helpful in measuring the solution degeneracy in feedforward networks.

K Detailed characterization of OOD generalization performance

In addition to showing the behavioral degeneracy in the main text, here we provide a more detailed characterization of the OOD behavior of networks by showing the mean versus standard deviation, and the distribution of the OOD losses.

K.1 Changing task complexity

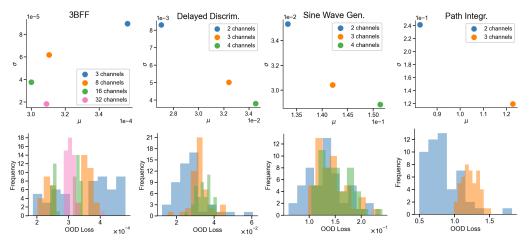


Figure 12: Detailed characterization of the OOD performance of networks while changing task complexity.

1 K.2 Changing feature learning strength

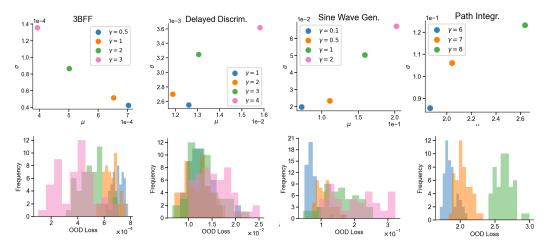


Figure 13: Detailed characterization of the OOD performance of networks while changing feature learning strength. Across Delayed Discrimination, Sine Wave Generation, and Path Integration tasks, networks trained with larger γ – and thus undergoing stronger feature learning – exhibit higher mean OOD generalization loss together with higher variability, potentially reflecting overfitting to the training task.

K.3 Changing network size

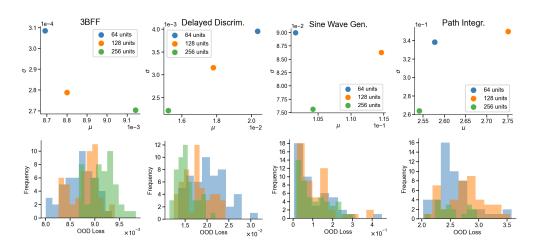


Figure 14: Detailed characterization of the OOD performance of networks while changing network size.

3 K.4 Changing regularization strength

K.4.1 Low-rank regularization

504

506

507

508

509

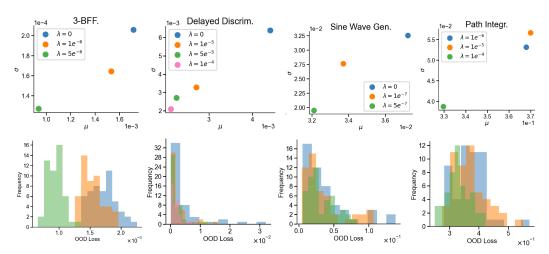


Figure 15: Detailed characterization of the OOD performance of networks while changing low-rank regularization strength.

5 K.4.2 Sparsity (L1) regularization

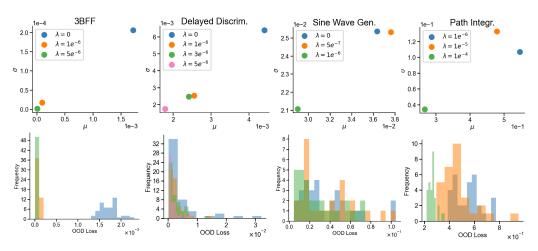


Figure 16: Detailed characterization of the OOD performance of networks while changing sparsity (L1) regularization strength.

L A short introduction to Maximal Update Parameterization (μP)

Under the NTK parametrization, as the network width goes to infinity, the network operates in the *lazy* regime, where its functional evolution is well-approximated by a first-order Taylor expansion around the initial parameters [49, 37, 32, 33]. In this limit feature learning is suppressed and training dynamics are governed by the fixed Neural Tangent Kernel (NTK).

dynamics are governed by the fixed Neural Tangent Kernel (NTK).

To preserve non-trivial feature learning at large width, the *Maximal Update Parametrization* (μ P)
rescales both the weight initialisation and the learning rate. μ P keeps three quantities *width-invariant*at every layer—(i) the norm/variance of activations (ii) the norm/variance of the gradients, and (iii)
the parameter updates applied by the optimizer [67, 68, 34, 35].

For recurrent neural networks, under Stochastic Gradient Descent (SGD), the network output, initialization, and learning rates are scaled as

$$f = \frac{1}{\gamma_0 N} \vec{w} \cdot \phi(h), \tag{1}$$

$$\partial_t h = -h + \frac{1}{\sqrt{N}} J \phi(h), \qquad J_{ij} \sim \mathcal{N}(0, 1),$$
 (2)

$$\eta_{\text{SGD}} = \eta_0 \, \gamma_0^2 \, N. \tag{3}$$

517 Under Adam optimizer, the network output, initialization, and learning rates are scaled as

$$f = \frac{1}{\gamma_0 N} \vec{w} \cdot \phi(h), \tag{4}$$

$$\partial_t h = -h + \frac{1}{N} J \phi(h), \qquad J_{ij} \sim \mathcal{N}(0, N), \tag{5}$$

$$\eta_{\text{Adam}} = \eta_0 \, \gamma_0. \tag{6}$$

518 M Theoretical relationship between parameterizations

- We compare two RNN formalisms used in different parts of the main manuscript: a standard discrete-
- time RNN trained with fixed learning rate and conventional initialization, and a μ P-style RNN trained
- with leaky integrator dynamics and width-aware scaling.
- In the standard discrete-time RNN, the hidden activations are updated as

$$h(t+1) = \phi(W_h h(t) + W_x x(t)),$$

In μP RNNs, the hidden activations are updated as

$$h(t+1) - h(t) = \tau \left(-h(t) + \frac{1}{N}J\phi(h(t)) + Ux(t)\right)$$

524 When $\tau = 1$,

525

528

529

530

531

$$h(t+1) - h(t) = -h(t) + \frac{1}{N}J\phi(h(t)) + Ux(t)$$

$$h(t+1) = \frac{1}{N}J\phi(h(t)) + Ux(t)$$

- Aside from the overall scaling factor, the difference between the two parameterizations lies in the placement of the non-linearity:
 - Standard RNN: φ is applied post-activation, i.e. after the recurrent and input terms are linearly combined,
 - μP RNN: φ is applied pre-activation; i.e. before the recurrent weight matrix, so the hidden state is first non-linearized and then linearly combined
- Miller and Fumarola [69] demonstrated that two classes of continuous-time firing-rate models which differ in their placement of the non-linearity are mathematically equivalent under a change of variables:

v-model
$$\tau \frac{dv}{dt} = -v + \tilde{I}(t) + W f(v)$$

r-model:
$$au \frac{dr}{dt} = -r + f(Wr + I(t))$$

- with equivalence holding under the transformation v(t) = Wr(t) + I(t) and $\tilde{I}(t) = I(t) + \tau \frac{dI}{dt}$, assuming matched initial conditions.
- Briefly, they show that Wr + I evolves according to the v-equation as follows:

$$\begin{split} v(t) &= Wr(t) + I(t) \\ \frac{dv}{dt} &= \frac{d}{dt} \left(Wr(t) + I(t) \right) \\ &= W \frac{dr}{dt} + \frac{dI}{dt} \\ &= W \left(\frac{1}{\tau} \left(-r + f(Wr + I) \right) \right) + \frac{dI}{dt} \\ \tau \frac{dv}{dt} &= -Wr + Wf(Wr + I) + \tau \frac{dI}{dt} \\ &= -(v - I) + Wf(v) + \tau \frac{dI}{dt} \\ &= -v + I + \tau \frac{dI}{dt} + Wf(v) \\ \tau \frac{dv}{dt} &= -v + \tilde{I}(t) + Wf(v) \end{split}$$

This mapping applies directly to RNNs viewed as continuous-time dynamical systems and helps relate v-type μ P-style RNNs to standard discrete-time RNNs. It suggests that the μ P RNN (in v-type form) and the standard RNN (in r-type form) can be treated as different parameterizations of the same underlying dynamical system when:

Initialization scales are matched

- The learning rate is scaled appropriately with γ
 - Output weight norms are adjusted according to width

In summary, while a theoretical equivalence exists, it is contingent on consistent scaling across all components of the model. In this manuscript, we use the standard discrete-time RNNs due to its practical relevance for task-driven modeling community, while switching to μP to isolate the effect of feature learning and network size. Additionally, we confirm that the feature learning and network size effects on degeneracy hold qualitatively the same in standard discrete-time RNNs, unless where altering network width induces unstable and lazier learning in larger networks (Figure Q and R).

N Verifying larger γ reliably induces stronger feature learning in μP

In μP parameterization, the parameter γ interpolates between lazy training and rich, feature-learning dynamics, without itself being the absolute magnitude of feature learning. Here, we assess feature-learning strength in RNNs under varying γ using two complementary metrics:

Weight-change norm which measures the magnitude of weight change throughout training. A larger weight change norm indicates that the network undergoes richer learning or more feature learning.

$$\frac{\|\mathbf{W}_T - \mathbf{W}_0\|_F}{N},$$

where N is the number of parameters in the weight matrices being compared.

Kernel alignment (KA), which measures the directional change of the neural tangent kernel (NTK) before and after training. A lower KA score corresponds to a larger NTK rotation and thus stronger feature learning.

$$\mathrm{KA}(K^{(f)}, K^{(0)}) = \frac{\mathrm{Tr}(K^{(f)}K^{(0)})}{\|K^{(f)}\|_{F} \|K^{(0)}\|_{F}}, \qquad K = \nabla_{W}\hat{y}^{\mathsf{T}}\nabla_{W}\hat{y}.$$

We demonstrate that higher γ indeed amplifies feature learning inside the network.

562 N.1 N-BFF

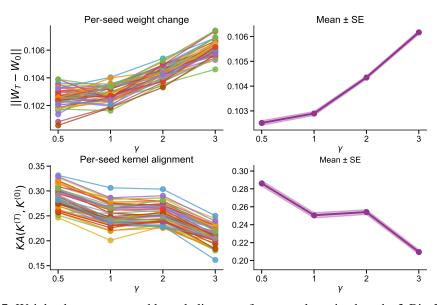


Figure 17: Weight change norm and kernel alignment for networks trained on the 3-Bits Flip Flop task as we vary γ . On the left panels, we show the per-seed metrics where connected dots of the same color are networks of identical initialization trained with different γ . On the right panels, we show the mean and standard error of the metrics across 50 networks. For larger γ , the weights move further from their initializations as shown by the larger weight change norm, and their NTK evolves more distinct from the network's NTK at initialization as shown by the reduced KA. Both indicate stronger feature learning for networks trained under larger γ .

N.2 Delayed Discrimination

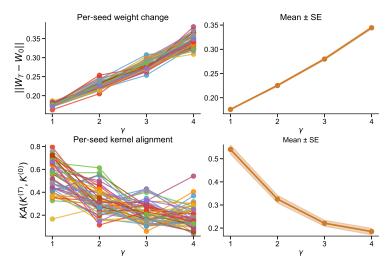


Figure 18: Stronger feature learning for networks trained under larger γ on the Delayed Discrimination task.

564 N.3 Sine Wave Generation

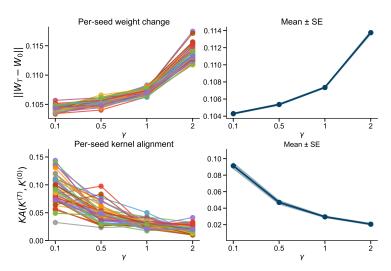


Figure 19: Stronger feature learning for networks trained under larger γ on the Sine Wave Generation task.

565 N.4 Path Integration

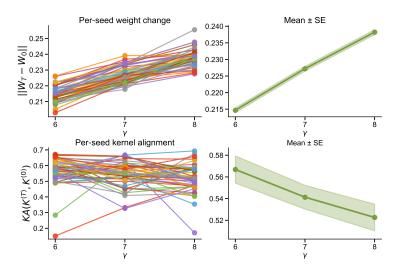


Figure 20: Stronger feature learning for networks trained under larger γ on the Path Integration task.

566 O Verifying μP reliably controls for feature learning across network width

Here, we only use Kernel Alignment to assess the feature learning strength in the networks since the unnormalized weight-change norm $\|\mathbf{W}_T - \mathbf{W}_0\|_F$ scales directly with matrix size (therefore network size) and there exists no obvious way to normalize across different dimensions. In our earlier analysis where we compared weight-change norms at varying γ , network size remained fixed, so those Frobenius-norm measures were directly comparable. We found that, for all tasks except Delayed Discrimination, the change in mean KA across different network sizes remains extremely small (less than 0.1), which demonstrates that μP parameterization with the same γ has effectively controlled for feature learning strength across network sizes. On Delayed Discrimination, the networks undergo slightly lazier learning for larger network sizes. Nevertheless, we still include Delayed Discrimination in our analyses of solution degeneracy to ensure our conclusions remain robust even when μP can't perfectly equalize feature-learning strength across widths. As shown in the main paper, lazier learning regime generally increases dynamical degeneracy; yet, larger networks which exhibit lazier learning in the N-BFF task actually display lower dynamical degeneracy. This reversed trend confirms that the changes in solution degeneracy arise from network size itself, not from residual variation in feature learning strength.

32 O.1 N-BFF

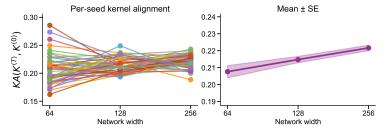


Figure 21: Kernel alignment (KA) for different network width on the 3 Bits Flip-Flop task. (Lower KA implies more feature learning.)

583 O.2 Delayed Discrimination

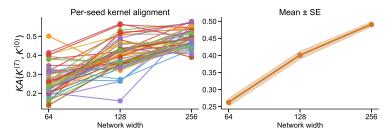


Figure 22: Kernel alignment for different network width on the Delayed Discrimination task.

584 O.3 Sine Wave Generation

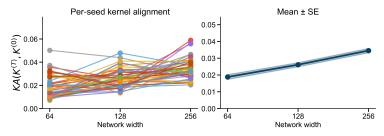


Figure 23: Kernel alignment for different network width on the Sine Wave Generation task.

585 O.4 Path Integration

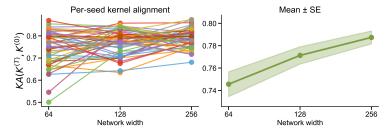


Figure 24: Kernel alignment for different network width on the Path Integration task.

P Regularization's effect on degeneracy for all tasks

In addition to showing regularization's effect on degeneracy in Delayed Discrimination task in the main paper, here we show that heavier low-rank regularization and sparsity regularization also reliably reduce solution degeneracy across neural dynamics, weights, and OOD behavior in the other three tasks.

591 P.1 Low-rank regularization

587

588

589

590

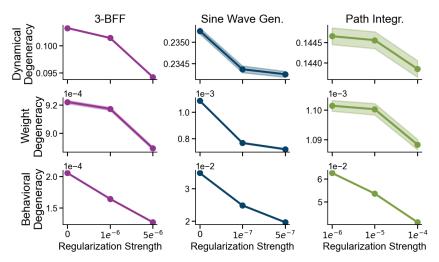


Figure 25: Low-rank regularization reduces degeneracy across neural dynamics, weight, and OOD behavior on the N-BFF, Sinewave Generation, and Path Integration task.

92 P.2 Sparsity regularization

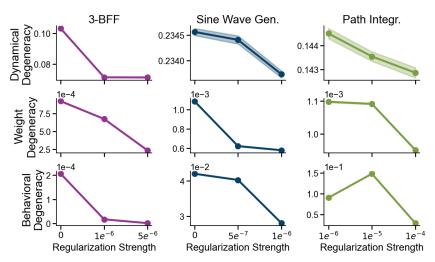


Figure 26: Sparsity regularization reduces degeneracy across neural dynamics, weight, and OOD behavior on the N-BFF, Sinewave Generation, and Path Integration task.

93 Q Test feature learning effect on degeneracy in standard parameterization

While μP lets us systematically vary feature-learning strength to study its impact on solution 594 degeneracy, we confirm that the same qualitative pattern appears in standard discrete-time RNNs: 595 stronger feature learning lowers dynamical degeneracy and raises weight degeneracy (Figure 27). 596 To manipulate feature-learning strength in these ordinary RNNs we applied the γ -trick—scaling the 597 network's outputs by γ —and multiplied the learning rate by the same factor. With width fixed, these 598 two operations replicate the effective changes induced by μP . Figure 28 shows that this combination 599 reliably tunes feature-learning strength. Besides weight-change norm and kernel alignment, we also 600 report representation alignment (RA), giving a more fine-grained view of how much the learned 601 features deviate from their initialization [39]. Representation alignment is the directional change of 602 the network's representational dissimilarity matrix before and after training, and is defined by

$$\mathrm{RA}\!\big(R^{(T)},R^{(0)}\big) := \frac{\mathrm{Tr}\!\big(R^{(T)}R^{(0)}\big)}{\|R^{(T)}\|\,\|R^{(0)}\|}, \qquad R \; := \; H^\top H,$$

A lower RA means more change in the network's representation of inputs before and after training, and indicates stronger feature learning.

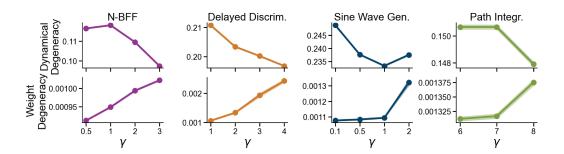


Figure 27: Stronger feature learning reliably decreases dynamical degeneracy while increasing weight degeneracy in standard discrete-time RNNs.

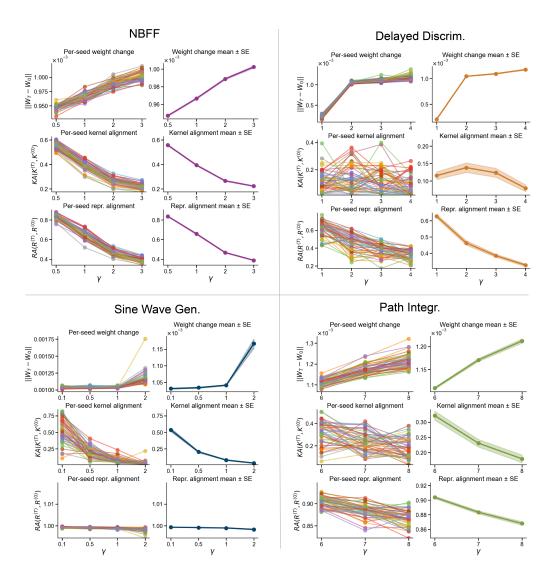


Figure 28: Larger γ reliably induces stronger feature learning in standard discrete-time RNNs.

R Test network size effect on degeneracy in standard parameterization

When we vary network width, both the standard parameterization and μP parameterization display the same overall pattern: larger networks exhibit lower dynamical and weight degeneracy. An exception arises in the 3BFF task, where feature learning becomes unstable and collapses in the wider models. In that setting we instead see *higher* dynamical degeneracy, which we suspect because the feature learning effect (lazier learning leads to higher dynamical degeneracy) dominates the network size effect.

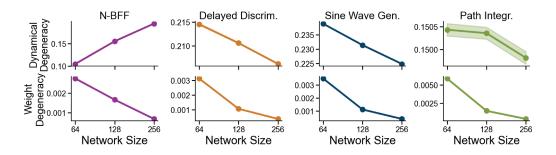


Figure 29: Larger network sizes lead to lower dynamical and weight degeneracy, except in the case where feature learning is unstable across width (in N-BFF).

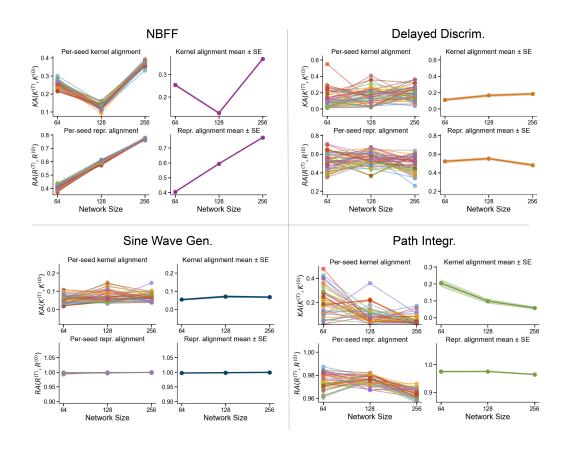


Figure 30: When changing network width in standard discrete-time RNNs, feature learning strength remains stable across width except in N-BFF, where notably lazier learning happens in the widest network.

S Disclosure of compute resources

In this study, we conducted 50 independent training runs on each of four tasks, systematically sweeping four factors that modulate solution degeneracy—task complexity (15 experiments), learning regime (15 experiments), network size (12 experiments), and regularization strength (26 experiments), resulting in a total of 3400 networks. Each experiment was allocated 5 NVIDIA V100/A100 GPUs, 32 CPU cores, 256 GB of RAM, and a 4-hour wall-clock limit, for a total compute cost of approximately 68 000 GPU-hours.