

## Abstract

1 By training linear physical networks to learn linear transformations, we discern  
2 how their physical properties evolve due to weight update rules. Our findings  
3 highlight a striking similarity between the learning behaviors of such networks  
4 and the processes of aging and memory formation in disordered and glassy sys-  
5 tems. We show that the learning dynamics resembles an aging process, where the  
6 system relaxes in response to repeated application of the feedback boundary forces  
7 in presence of an input force, thus encoding a memory of the input-output rela-  
8 tionship. With this relaxation comes an increase in the correlation length, which is  
9 indicated by the two-point correlation function for the components of the network.  
10 We also observe that the square root of the mean-squared error as a function of  
11 epoch takes on a non-exponential form, which is a typical feature of glassy sys-  
12 tems. This physical interpretation suggests that by encoding more detailed infor-  
13 mation into input and feedback boundary forces, the process of emergent learning  
14 can be rather ubiquitous and, thus, serve as a very early physical mechanism, from  
15 an evolutionary standpoint, for learning in biological systems.

## 16 Emergent learning in physical systems as feedback-based aging in a glassy 17 landscape

18 Vidyesh Rao Aniseti,<sup>1</sup> Ananth Kandala,<sup>2</sup> J. M. Schwarz<sup>1,3</sup>

19 <sup>1</sup>*Physics Department, Syracuse University, Syracuse, NY 13244 USA*

20 <sup>2</sup>*Department of Physics, University of Florida, FL 32611-8440, USA*

21 <sup>3</sup>*Indian Creek Farm, Ithaca, NY 14850 USA*

22 (Dated: October 4, 2023)

## 23 I. INTRODUCTION

24 Given the prevalence of emergent behavior, physicists, computer scientists, and biologists have long  
25 asked whether or not some subset of emergent behavior results in the capacity of a system of many  
26 interacting components to learn, i.e., to have intelligence [1, 2]. While there has been much focus  
27 looking for emergent learning in brain-like systems, such as neuronal networks in biology or artifi-  
28 cial neural networks in physics and computer science, recent research has demonstrated that simple  
29 physical systems, such as a spring network, have the potential to exhibit learning behavior similar  
30 to that of artificial neural networks [3–9]. In this context, learning refers to the ability to modify  
31 the properties of a physical system by adjusting its learning degrees of freedom in order to more  
32 efficiently achieve some task. For example, in a spring network, the spring stiffness and rest lengths  
33 represent the learning degrees of freedom, while the nodes of the springs correspond to the usual  
34 physical degrees of freedom.

35 In these physical learning systems, once input boundary nodes, output boundary nodes, and a cost  
36 function are all chosen, the learning process is composed of two steps; 1) *Signaling* : System’s  
37 response to a given input is compared with the desired output and an update signal is sent which  
38 provides information on the necessary adjustments to each learning degree of freedom, so that the  
39 system’s response aligns more closely with the desired output. 2) *Weight update* : Each learning  
40 degree of freedom, or weight, is updated in response to the update signal. This weight update should  
41 allow the system to perform gradient descent. The two steps are repeatedly applied to train the  
42 system to learn.

43 The major challenge applying this algorithm is to find physical processes that implement the above  
44 two steps. While methods such as Equilibrium Propagation (EP) [4], Multi-mechanism Learning  
45 (MmL) [3, 5], and Coupled Learning (CP) [6] have made strides in addressing this challenge, they  
46 are not entirely physical in nature. In particular, the learning stages involved, *Signaling* and *Weight*  
47 *update*, require the artificial modifications to the physical system. For instance, in EP and CL,  
48 to send the gradient information into the system, one needs to store the free state in some memory,  
49 which is not possible in typical systems such as spring networks or resistor networks. In our previous

50 work unveiling MmL, we demonstrated that this issue of memory storage could be addressed by  
 51 encoding the feedforward and feedback signal into two non-interfering physical quantities [3, 5].

52 Despite this demonstration, however, a significant problem remains: we do not know of any physical  
 53 process that can update the weights in the system. To physically implement weight updates, recent  
 54 experimental efforts have resorted to using complex components such as transistors in the training  
 55 of electrical networks [10, 11], and actuators and sensors in mechanical networks [12]. Yet, the  
 56 reliance on such intricate and varied tools introduces challenges in terms of scalability and robust-  
 57 ness in these approaches. Here, we explore the central question: Do effects of the weight update  
 58 procedure resemble any natural physical phenomena? The answer to such a question will point us  
 59 in the direction of a fully physical learning system, weight update included. To begin to answer  
 60 this question, we train linear physical networks and investigate how the physical properties of this  
 61 system change, given the weight update rule.

62 Our manuscript consists of revisiting our MmL training procedure, as detailed in our prior work [3,  
 63 5], in a general manner that emphasizes its physical plausibility in Section II. Results are then  
 64 presented in Section III. We conclude with a discussion of the impact of our results. ( We also  
 65 review the specifics of multi-mechanism learning in Appendix A, followed by data generation and  
 66 network generation in Appendix B.)

## 67 II. THE LEARNING PROCESS

68 We now demonstrate the process of physical learning within our system. Initially, we impose an  
 69 input boundary condition, denoted by  $I$ . The system's response is then captured by the Laplace  
 70 equation  $Lv = I$ , where  $L$  is Laplacian, which depends on the learning degrees of freedom  $w$ , and  
 71  $v$  is the state of the system. To attain its intended functionality, the system need to update  $w$  to  
 72 minimize the cost function  $C(v(w))$ . We encode the cost function as an interaction energy between  
 73 the system and the environment. This energy causes a feedback boundary condition of the form  
 74  $-\eta \frac{\partial C(v)}{\partial v}$  to act on the system, due to which the state of the system evolves along a direction that  
 75 decreases  $C(v)$ :

$$L(v + \delta v) = I - \eta \frac{\partial C(v)}{\partial v}. \quad (1)$$

76 For a mechanical network, these input and feedback boundary conditions are applied as external  
 77 stresses on the system. When the feedback stress is removed, the system tends to revert to its initial  
 78 state  $v$ . However, with continuous exposure to feedback boundary forces, there's a lasting change in  
 79 the system's learning degrees of freedom. This change is akin to a plastic deformation in materials  
 80 where repeated stress leads to permanent alterations.

81 Note that unlike the input boundary condition, the feedback boundary condition is a function of the  
 82 state of the system. As a result, there exists an optimal state where the system experiences minimal  
 83 feedback stress. Our hypothesis is that, through repeated application of these feedback stresses, the  
 84 system's learning parameters  $w$  evolve such that this optimal state is reached. The objective of this  
 85 evolution is to minimize the external stress  $-\eta \frac{\partial C(v)}{\partial v}$ , by changes in state of the system  $v$ , through  
 86 changes in  $w$ . This adaptation is represented as:

$$\Delta w_{ij} = -\alpha \eta \frac{\partial C(w)}{\partial w_{ij}}, \quad (2)$$

87 where  $C$  is a function of  $w$  via  $C(v(w))$ . In our previous work [3], we showed that the above weight  
 88 update rule can be written purely in terms of local physical quantities

$$\Delta w_{ij} = -\alpha v_{ij} \delta v_{ij}. \quad (3)$$

89 Where,  $w_{ij}$  is the weight connecting nodes  $(i, j)$ , and  $v_{ij}$  is the potential drop  $v_i - v_j$ ,  $\delta v_{ij}$  is  
 90 the change in this potential drop due to feedback[13]. Intriguingly, this learning rule exhibits a  
 91 Hebbian-like behavior.

92 Due to the evolution of the learning degrees of freedom, once reaching steady state, the sys-  
 94 tem's response is :

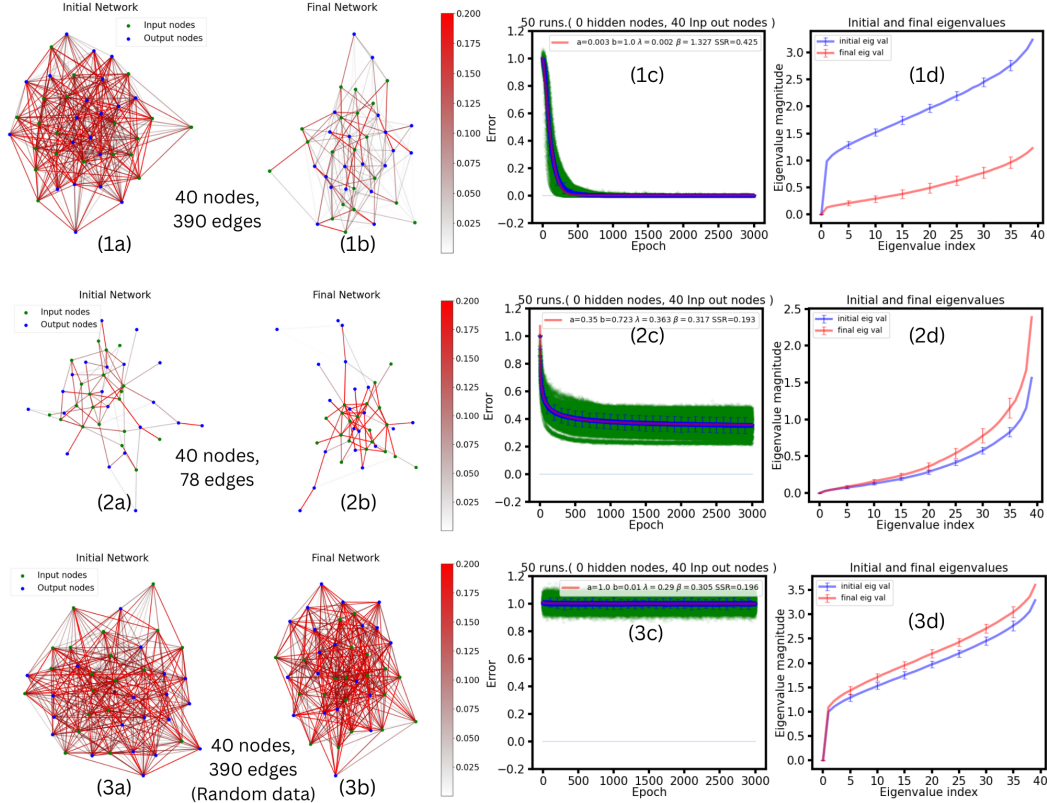


FIG. 1. *Training linear networks to learn linear transformations.* [1a & 1b] : *Network undergoes trimming.* A network with 40 nodes and 390 edges is trained to learn a linear transformation of size  $10 \times 10$ . Weights of the network are uniformly sampled from  $[10^{-5}, 0.2]$ . Colorbar on right shows weight values of each edge. [1c] *Non-exponential relaxation* : Training curve for the case shown in 1a and 1b but for 50 different initializations (shown in green). Y axis shows error defined as square root of mean square error, X axis shows epoch. In one epoch the network goes through 100 data points . All green curves are obtained after normalization with their respective initial errors. The blue curve shows the average over these 50 runs. The blue curve is fit to a non-exponential curve of the form  $a + be^{-\lambda \cdot t^\beta}$ . Fit parameters are shown in the legend.  $\beta > 1$  shows the relaxation shows a compressed exponential behaviour. The sum of squared residuals (SSR) is used to assess the goodness of fit, it is defined as:  $SSR = \sum_{i=1}^n (y_i^{fit} - y_i^{data})^2$  [1d] *Eigenvalues decrease while learning*: Eigenvalues of graph Laplacian before and after training for runs shown in 1c. These initial and final eigenvalues are averaged over those 50 runs. The eigenvalues are sorted in increasing order. The x-axis shows eigenvalue index. The network has 40 nodes so there are 40 eigenvalues. [2a, to d] These plots show the training performance for a network with less number of edges (78 edges), due to which it does not learn well. When compared with case 1, we see that trimming is less prominent and the eigenvalues do not decrease. The training curve shows a stretched exponential relaxation ( $\beta < 1$ ) and saturates well above zero error. [3a, to d] *Training on random data*: Networks initialized with same parameters as that of 1a are trained on randomly generated data. No trimming is observed, eigenvalues increase over training and the error curve does not decrease with the number of epochs.

$$L'(v + \delta v) = I, \quad (4)$$

95 where  $L'$  is the updated Laplacian that encodes the memory of the feedback stress by adapting to it,  
 96 i.e;  $C(v + \delta v) < C(v)$ .

97 In summary, the learning process goes as follows. An input is introduced to the system as an external  
 98 force. Subsequently, based on the system's reaction to this input, feedback forces are consistently  
 99 applied. We postulate that such a process enables the system to adapt and become attuned to these  
 100 feedback boundary forces. This continuous adaptation to feedback forces, in presence of the input,  
 101 ingrains a memory of the input-output relationship within the system. This concept is elucidated  
 102 further in the subsequent section.

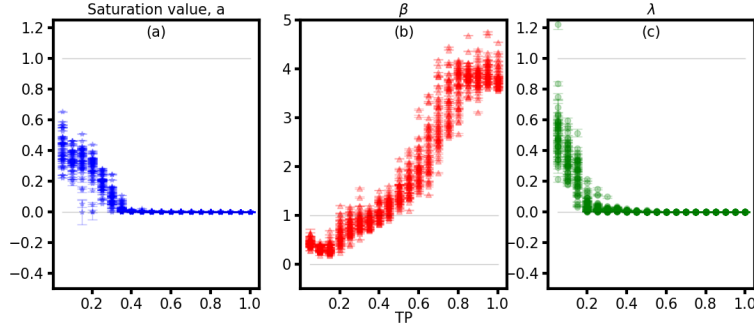


FIG. 2. *Learning performance with overparametrization.* Error curve is fit to  $a + be^{-\lambda t^\beta}$  for networks with varying edges and the fit parameters are plotted (Error bars shown are calculated using the diagonal terms of covariance matrix). The Tuning Parameter ( $TP$ ) serves as a metric to quantify the degree of connectivity in a network. Specifically, it is calculated by taking the ratio of the number of edges  $M$  present in the graph to the number of edges that would exist in a fully connected network with the same number of nodes. (a) We observe that after adding a certain number of edges, the saturation value of the error curve begins to asymptote to zero. (b) We also observe that the exponent  $\beta$  increases from less than one to greater than one, showing a shift from stretched exponential to compressed exponential relaxation. (c)  $\lambda$  value also becomes very small after adding a certain number of edges. We have done a fit robustness analysis for these plots in Appendix A. (In Fig.1, 390 and 78 edge networks correspond to a  $TP$  of 0.5 and 0.1, respectively.)

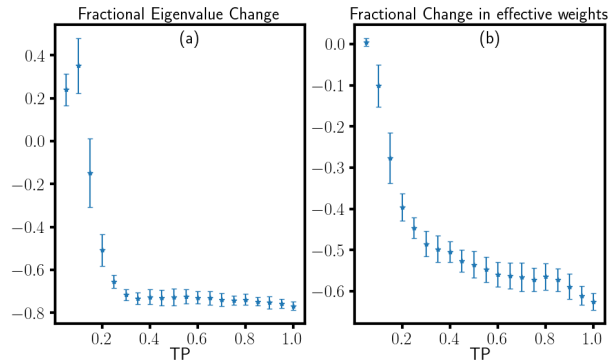


FIG. 3. *Eigenvalue decrease and trimming with overparametrization.* (a) Shows fractional decrease in the sum of eigenvalues due to learning, averaged over 50 runs. (b) Shows fractional decrease in number of effective weights due to learning, averaged over 50 runs. Here, the term ‘effective weights’ refers to those weights that fall within the top 99 percent of the permissible weight value range( $[10^{-5}, 0.2]$ ).

### 103 III. RESULTS

104 Figure 1.1a,1b shows the network before and after training for a network of  $N = 40$  and  $M = 390$ .  
 105 Since the intensity of the color indicates the magnitude of the weight, note that many of the weights  
 106 of a trained network reach the minimum value. In other words, there is a trimming effect, where  
 107 only the important edges remain. To ascertain whether or not the network has learned the linear  
 108 transformation, we plot the square root of the mean-squared error in Fig. 1.1c as a function of epoch.  
 109 Given that the error nearly vanishes at longer epoch, this network has successfully learned the task.  
 110 This shows the dynamics through which the system relaxes to the feedback boundary forces due  
 111 to the evolution of learning degrees of freedom. Interestingly, we performed a phenomenological  
 112 fit for this curve. The curve is well-approximated by a non-exponential relaxation of the form  
 113  $\sqrt{MSE} = a + b \exp(-\lambda t^\beta)$ , where  $a, b, \lambda, \beta$  are the fit parameters and  $t$  denotes the epoch  
 114 number. Interestingly, these dynamics are quantitatively similar to what is observed in molecular  
 115 glassy systems [14]. This finding demonstrates the existence of a glassy landscape. Appendix A  
 116 addresses the reasonableness of this non-exponential fit.

117 We seek to quantify further the relaxation of the system as it learns. We, therefore, compute the  
 118 eigenvalues of the Laplacian matrix. Figure 1.1d shows how learning results in decreasing Laplacian

119 eigenvalues. Note that these eigenvalues are the square root of normal mode frequencies. Decreasing  
 120 eigenvalues is evidence that the network is getting “softer” as the normal mode excitations become  
 121 longer in wavelength. This observation demonstrates that the network moves from a state of stress  
 122 to that of less stress due to repeated application feedback boundary forces. The network is, thus,  
 123 “adapting” to these feedback forces indicating a transition towards a state that encodes a memory  
 124 of the input-output relationship. Additionally, it draws parallels between this behavior and the self-  
 125 organization observed in periodically sheared suspensions, where the system adapts to the periodic  
 126 driving in a similar manner [15]. Moreover, when amorphous solids, modeled as purely repulsive  
 127 particles in the jammed phase, are shear-stabilized by minimizing the energy with respect to the  
 128 shear degrees of freedom, one finds longer wavelength excitations emerging [16]. Finally, recent  
 129 work demonstrates that using a similar multiplicative learning rule as given in Eq. 7 to train physical  
 130 networks to learn linear transformations also shows a decrease in the lowest eigenvalues of the  
 131 Hessian [17]. Appendix D shows that the trends hold for larger system sizes.

132 Figures 1.2(a-d) show the same quantities as Figure 1.1, however, for a network with  $N = 40$  and  
 133  $M = 78$ . Given the smaller number of learning degrees of freedom, a network with this architecture  
 134 does not successfully learn, as indicated by the square root of the mean-squared error not decreasing  
 135 to zero as the number of epochs increase. Moreover, the eigenvalues of the Laplacian do not decrease  
 136 and so the system does not relax, or soften. For comparative purposes, we also train the network  
 137 to learn, if you will, random data. Fig. 1.(3a to d) shows the physical effects of learning random  
 138 data. Here, the system, exposed to random input and feedback boundary conditions, does not relax,  
 139 as indicated by the unchanged initial and final eigenvalues. With random input-output forces, the  
 140 weight update signal in Eq. 3 averages to zero due to the absence of correlation between  $v_{ij}$  and  
 141  $\delta v_{ij}$ . This null result suggests that the system’s relaxation is driven by correlations between input  
 142 and feedback boundary conditions and for certain network architectures.

143 Given the nontrivial dependence of learning on the network architecture, we further extend our  
 144 analysis by incrementally increasing the network connectivity to examine the implications of over-  
 145 parametrization (see Fig. 2). We denote the ratio of the number of edges  $M$  to the number of nodes  
 146 in the fully connected equivalent network as  $TP$  for tuning parameter. The results indicate that  
 147 as more edges are introduced, the cost landscape becomes steeper due to a reduced number of flat  
 148 directions [18], leading to accelerated relaxation and enhanced learning performance. Notably, a  
 149 parallel can be drawn with glasses; in these systems, increased connectivity also speeds up relax-  
 150 ation dynamics [19, 20]. Both these studies, as well as ours, show a shift in relaxation dynamics  
 151 from a stretched to a compressed exponential upon increasing connectivity. This further underscores  
 152 the intrinsic link between learning processes and relaxation in disordered systems.

153 Given the changes in the weights as the networks learns, in Fig. 3, we examine the relationship be-  
 154 tween trimming, eigenvalue reduction, and network connectivity. As network connectivity increases  
 155 by increasing  $TP$ , the fractional eigenvalue decrease tends to plateau, reaching a saturation point  
 156 around  $TP \approx 0.3$ . A comparison of Fig. 3(a) and Fig. 2(a) reveals a notable correlation: the point of  
 157 eigenvalue saturation aligns with the disappearance of saturation error. This suggests a fundamental  
 158 link between the processes of learning and eigenvalue reduction. Furthermore, Fig. 3(b) underscores  
 159 the ubiquity of the trimming effect across networks of varying connectivity. Notably, the magnitude  
 160 of the trimming effect intensifies as network size grows.

161 Figure 4 illustrates the evolution of the resistance distance distribution during the learning process.  
 162 In an electrical network, the effective resistance between two nodes can be interpreted as a measure  
 163 of distance (more details in Appendix C). By calculating the average distribution of resistance dis-  
 164 tances over all possible pairs of nodes, a two-point correlation function  $p(r)$  can be derived, which  
 165 can be extended to spring and flow networks as well. As learning progresses, we observe a broad-  
 166 ening of the two-point correlation function, indicating that the average conductance between two  
 167 arbitrary nodes decreases. This phenomenon is analogous to a reduction in “stiffness” in elastic  
 168 networks, as the system becomes more soft during learning.

## 169 IV. DISCUSSION

170 In summary, in learning about the physical signatures of multi-mechanism learning we find that; 1)  
 171 The error curve for networks with low connectivity resembles a stretched exponential. However, as  
 172 network connectivity increases, the error curve transitions to a compressed exponential form (Fig.

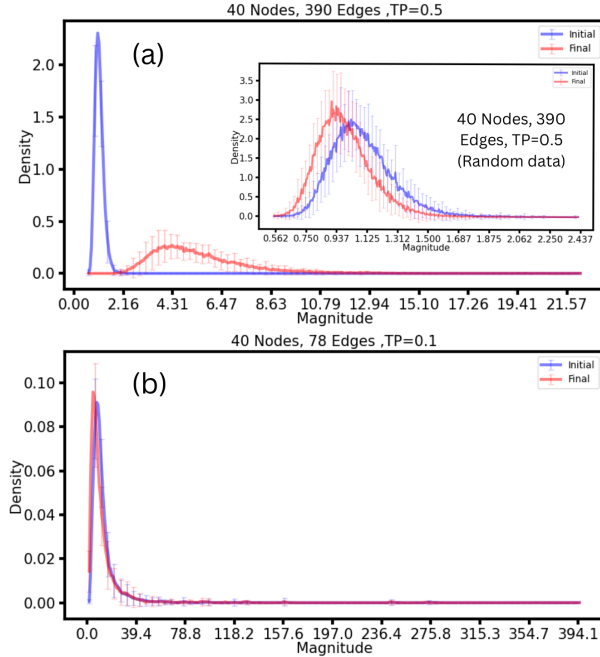


FIG. 4. *Resistance Distance Distribution and Learning*. (a) The figure showcases the average resistance distance distribution,  $p(r)$ , during learning, with the x-axis denoting resistance magnitude and the y-axis its normalized frequency. This is averaged over 50 network initializations. The inset illustrates the outcome when the network is trained on random data (note that the scale in the inset differs, making the initial distributions appear distinct, though they are identical). (b) Represents a network with suboptimal learning performance due to a limited number of edges.

173 2). 2) Eigenvalues of the graph Laplacian decrease with epoch and long wavelength modes are  
 174 generated (Fig. 1). 3) The network undergoes trimming, i.e., lot of the weights go to zero (Fig. 1 &  
 175 3). 4) The two point correlation function for the network broadens while learning (Fig. 4).

176

177 Neuromorphic researchers have been actively seeking physical counterparts to facilitate au-  
 178 tonomous weight updates. This pursuit has led to the development of physical learning systems  
 179 utilizing memristors [21], nanoscale devices [22], and transistors [23]. However, the intricate  
 180 design requirements for each component presents challenges in terms of robustness and scalability.  
 181 We propose that soft materials might offer a more streamlined solution. These materials inherently  
 182 exhibit self-adjustment to external conditions, as evidenced by the self-organization of granular  
 183 systems in response to external driving [15, 24, 25] the adaptability of other disordered systems to  
 184 external strain [26, 27]. Consequently, they emerge as promising candidates for crafting physical  
 185 learning systems. Moreover, the model introduced in Section II provides insights into a potential  
 186 training methodology for soft materials, be it particulate-based, such as a granular learner, where  
 187 the topology of the system can change, or spring-based, such a spring network learner, where the  
 188 topology of the network is fixed. By iteratively applying input and feedback boundary forces,  
 189 the learning parameters can autonomously adapt to these forces to optimize a cost function. This  
 190 approach paves the way for the creation of innovative disordered materials with neural network-like  
 191 learning potential. We aim to validate this concept in our forthcoming research.

192 Finally, by using multi-mechanism learning to train physical networks to learn linear transforma-  
 193 tions, we demonstrate a simple, brain-like task in a typically non-brain-like material. As brains  
 194 began to emerge several hundred million years ago in planarians [28], physical learning mecha-  
 195 nisms are ripe candidates for life learning to survive in their environment before planarians. We,  
 196 therefore, seek to validate such mechanisms in pre-planarian organisms.

197 The authors thank Benjamin Scellier, Arvind Murugan, Eli Hawkins, Shabeeb Ameen and Samuel  
 198 Ropert for helpful discussion. JMS acknowledges financial support from NSF-DMR-2204312.

- 
- 199 [1] J. J. Hopfield, “Neural networks and physical systems with emergent collective computational abilities.,”  
 200 *Proceedings of the national academy of sciences*, vol. 79, no. 8, pp. 2554–2558, 1982.
- 201 [2] W. D. Hillis, “Intelligence as an emergent behavior; or, the songs of eden,” *Daedalus*, pp. 175–189, 1988.
- 202 [3] V. R. Anisetti, B. Scellier, and J. M. Schwarz, “Learning by non-interfering feedback chemical signaling  
 203 in physical networks,” *Phys. Rev. Research*, vol. 5, p. 023024, 2023.
- 204 [4] B. Scellier, *A deep learning theory for neural networks grounded in physics*. PhD thesis, Université de  
 205 Montréal, 2021.
- 206 [5] V. R. Anisetti, A. Kandala, B. Scellier, and J. M. Schwarz, “Frequency propagation: Multi-mechanism  
 207 learning in nonlinear physical networks,” *arXiv preprint arXiv:2208.08862*, 2022.
- 208 [6] M. Stern, W. Bialek, J. W. Shaevitz, M. Pan, H. Zhang, and A. Murugan, “Supervised learning in physical  
 209 networks: From machine learning to learning machines,” *arXiv preprint arXiv:1804.10130*, 2021.
- 210 [7] J. Kendall, R. Pantone, K. Manickavasagam, Y. Bengio, and B. Scellier, “Training end-to-end analog  
 211 neural networks with equilibrium propagation,” *arXiv preprint arXiv:2006.01981*, 2020.
- 212 [8] M. Stern and A. Murugan, “Learning without neurons in physical systems,” 2022.
- 213 [9] M. Stern, C. Arinze, L. Perez, S. E. Palmer, and A. Murugan, “Supervised learning through physical  
 214 changes in a mechanical system,” *Proceedings of the National Academy of Sciences*, vol. 117, no. 26,  
 215 pp. 14843–14850, 2020.
- 216 [10] S. Dillavou, B. Beyer, M. Stern, M. Z. Miskin, A. J. Liu, and D. J. Durian, “Circuits that train themselves:  
 217 decentralized, physics-driven learning,” in *AI and Optical Data Sciences IV*, vol. 12438, p. 124380G,  
 218 SPIE OPTO, 2023.
- 219 [11] S. Dillavou, M. Stern, A. J. Liu, and D. J. Durian, “Demonstration of decentralized physics-driven learn-  
 220 ing,” *Phys. Rev. Applied*, vol. 18, p. 014040, 2022.
- 221 [12] R. H. Lee, E. A. B. Mulder, and J. B. Hopkins, “Mechanical neural networks: Architected materials that  
 222 learn behaviors,” *Science Robotics*, 2022.
- 223 [13] We could also have defined  $v_{ij} = v_j - v_i$ , note that the learning rule is independent of this choice.
- 224 [14] J. C. Phillips, “Stretched exponential relaxation in molecular and electronic glasses,” *Reports on Progress  
 225 in Physics*, vol. 59, p. 1133, 1996.
- 226 [15] L. Corté, P. M. Chaikin, J. P. Gollub, and D. J. Pine, “Random organization in periodically driven systems,”  
 227 *Nature Physics*, vol. 4, pp. 420–424, 2008.
- 228 [16] H. Mizuno, H. Shiba, and A. Ikeda, “Continuum limit of the vibrational properties of amorphous solids,”  
 229 *Proceedings of the National Academy of Sciences*, vol. 114, no. 46, pp. E9767–E9774, 2017.
- 230 [17] M. Stern, A. J. Liu, and V. Balasubramanian, “The physical effects of learning,” *bioRxiv*, pp. 2023–06,  
 231 2023.
- 232 [18] M. Baity-Jesi, L. Sagun, M. Geiger, S. Spigler, G. B. Arous, C. Cammarota, Y. LeCun, M. Wyart, and  
 233 G. Biroli, “Comparing dynamics: deep neural networks versus glassy systems,” *Journal of Statistical  
 234 Mechanics: Theory and Experiment*, vol. 2019, p. 124013, December 2019.
- 235 [19] B. Cui, R. Milkus, and A. Zaccone, “The relation between stretched-exponential relaxation and the vi-  
 236 brational density of states in glassy disordered systems,” *Physics Letters A*, vol. 381, no. 4, pp. 338–343,  
 237 2016.
- 238 [20] Z. W. Wu, W. Kob, W.-H. Wang, and L. Xu, “Stretched and compressed exponentials in the relaxation  
 239 dynamics of a metallic glass-forming melt,” *Nature Communications*, vol. 9, 2018.
- 240 [21] D. B. Strukov, G. S. Snider, D. R. Stewart, and R. S. Williams, “The missing memristor found,” *Nature*,  
 241 vol. 453, pp. 80–83, 2008.
- 242 [22] D. Kuzum, R. G. D. Jeyasingh, B. Lee, and H.-S. P. Wong, “Nanoelectronic programmable synapses  
 243 based on phase change materials for brain-inspired computing,” *Nano Letters*, vol. 12, no. 5, pp. 2179–  
 244 2186, 2012.
- 245 [23] M. Jerry, P. Y. Chen, J. Zhang, P. Sharma, K. Ni, S. Yu, and S. Datta, “Ferroelectric fet analog synapse  
 246 for acceleration of deep neural network training,” in *2017 IEEE International Electron Devices Meeting  
 247 (IEDM)*, pp. 6.2.1–6.2.4, IEEE, 2017.
- 248 [24] N. C. Keim and S. R. Nagel, “Generic transient memory formation in disordered systems with noise,”  
 249 *Phys. Rev. Lett.*, vol. 107, p. 010603, 2011.
- 250 [25] D. Hexner, A. J. Liu, and S. R. Nagel, “Periodic training of creeping solids,” *Physical Review E*, vol. 101,  
 251 no. 3, p. 032906, 2020.
- 252 [26] N. Pashine, D. Hexner, A. J. Liu, and S. R. Nagel, “Directed aging, memory, and nature’s greed,” *Science  
 253 Advances*, vol. 5, p. eaax4215, 2019.
- 254 [27] D. Hexner, N. Pashine, A. J. Liu, and S. R. Nagel, “Effect of directed aging on nonlinear elasticity and  
 255 memory formation in a material,” *Physical Review Research*, vol. 2, p. 043231, 2020.

- 256 [28] H. B. Sarnat and M. G. Netsky, “The brain of the planarian as the ancestor of the human brain,” *Canadian*  
 257 *Journal of Neurological Sciences*, vol. 12, no. 4, pp. 296–302, 1985.  
 258 [29] D. A. Spielman, *Spectral and Algebraic Graph Theory*, ch. 12.8. 2023.

## 259 Appendix A: A brief review of Multi-mechanism Learning

260 We study a network comprised of nodes and connected by weighted edges. Let us represent the  
 261 weight of the edge between node  $x$  and node  $y$  as  $w_{xy}$ , which could signify conductances in an  
 262 electrical network, spring constants in a mechanical spring network, or pipe thickness in a flow  
 263 network, etc.

264 **Input Nodes:** An “input” node pair is pair of nodes  $(b_j^+, b_j^-)$  such that an input current  $I_j$  enters the  
 265 network via node  $b_j^+$  and exits through  $b_j^-$ . (For mechanical networks input current can be thought  
 266 of as external forces acting at input nodes). Let there be  $q$  such input node pairs in the network,  
 267 denoted by  $\{(b_1^+, b_1^-), (b_2^+, b_2^-), \dots, (b_q^+, b_q^-)\}$ .

268 **Output Nodes:** In response to the input currents, the system develops an electric potential at each  
 269 node. The network’s output is defined to be the set of potential differences across certain “output”  
 270 node pairs, obtained as  $v(o_i^+, o_i^-) = v(o_i^+) - v(o_i^-)$  for each output node pair  $(o_i^+, o_i^-)$ . Let there  
 271 be  $p$  such output node pairs in the network, represented as  $\{(o_1^+, o_1^-), (o_2^+, o_2^-), \dots, (o_p^+, o_p^-)\}$ .

272 **Cost Function:** The goal of training is to adjust the weights  $\{w_{xy}\}$  so that for a given set of input  
 273 currents, the desired potential drops  $\{v_d(o_i^+, o_i^-)\}$  are achieved across all the output nodes. We  
 274 employ a Mean Squared Error (MSE) cost function:

$$C = \frac{1}{2} \sum_{i=1}^p (v(o_i^+, o_i^-) - v_d(o_i^+, o_i^-))^2. \quad (A1)$$

275 **Feedback Mechanism:** To optimize this cost function, we introduce a feedback signal into the  
 276 network at the output nodes. For each output node pair, the feed-back current is calculated as:

$$\epsilon_i = -\eta(v(o_i^+, o_i^-) - v_d(o_i^+, o_i^-)) \quad (A2)$$

277 This current enters the network through node  $o_i^+$  and exits via  $o_i^-$ , with  $\eta$  being a positive “nudging”  
 278 factor. The feedback currents change the potentials at each node and let the change in the potential  
 279 at node  $j$  be denoted by  $u_j$ .

280 **Weight Update Rule:** The weights are then updated as:

$$\Delta w_{xy} = -\alpha u(x, y) v(x, y), \quad (A3)$$

281 where  $\alpha$  is the learning rate. This rule effectively performs gradient descent on the cost function:

$$\Delta w_{xy} = -\alpha \eta \frac{\partial C}{\partial w_{xy}} \quad (A4)$$

282 **Considerations:** The weight update is local, and its sign depends on the potential drops due to  
 283 input and feedback. We assume the system’s relaxation time is much shorter than the weight update  
 284 time, ensuring a steady state during weight adjustments. The two quantities in the weight update  
 285 must be independent. This can be ensured by encoding them into distinct physical quantities[5].  
 286 (Further details on the learning procedure and its physical implementation are given in Ref.[3]). For  
 287 larger networks, a higher learning rate is necessary to maintain the magnitude of weight changes.  
 288 To address this, we conduct a trial run for one epoch, adjusting the learning rate to ensure  $|\Delta w| \approx$   
 289  $10^{-3}$ . Additionally, we impose regularization by (1) Limiting each weight update:  $|\Delta w_{xy}| < \epsilon$ , and  
 290 (2) Constraining weight values:  $w_{min} \leq w_{xy} \leq w_{max}$ . This ensures a smooth training process and  
 291 prevents weights from becoming too large or too small. In our simulations, we set  $w_{min} = 0.00001$ ,  
 292  $w_{max} = 0.2$ , and  $\epsilon = 0.01$ .

## 293 Appendix B: Methodology

294 **Network Generation:** We aim to create networks consisting of  $N$  nodes, with a varying number  
 295 of edges  $M$ . For this, we first create a Barabási-Albert network with connection parameter 1. This  
 296 graph generation algorithm connects a new node with 1 existing node in a manner that nodes with  
 297 higher degree have a stronger likelihood for selection. This creates a network with  $N$  nodes and



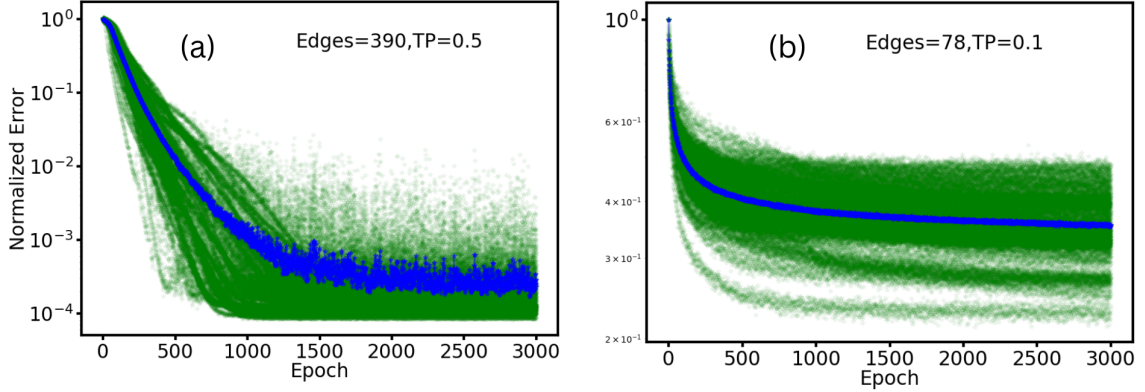


FIG. A1. *Log-Linear Analysis of Error Curves*: Panels (a) and (b) represent the log-linear plots corresponding to the error curves from Fig. 1: 2c and 1c, respectively. These analyses pertain to a 40-node network with 78 and 390 edges, respectively.

298  $N - 1$  edges. To create a network with  $M$  edges, we add  $M - (N + 1)$  unique edges. This way,  
 299 we can create networks with varying connectivity, ranging from being minimally connected to being  
 300 maximally connected. Note that it is highly unlikely to create such minimally connected networks  
 301 using the Erdős–Rényi model.

302 The generated networks are then trained on data generated using a linear transformation. Note that  
 303 in spite of using linear networks to learn linear transformations, the optimization needs to take place  
 304 in a cost landscape which is non-convex, high- dimensional, and disordered.

305 **Data Generation:** The input vector  $\mathbf{x}$  (eg;  $(x_1, x_2, x_3)$ ) is encoded as external currents across  
 306 input nodes  $\{(b_1^+, b_1^-), (b_2^+, b_2^-), (b_3^+, b_3^-)\}$  with currents  $+x_q$  and  $-x_q$  applied across nodes  $b_q^+$   
 307 and  $b_q^-$  respectively. The output vector  $\mathbf{y}$  (eg;  $(y_1, y_2, y_3)$ ) is the potential drop across nodes  
 308  $\{(o_1^+, o_1^-), (o_2^+, o_2^-), (o_3^+, o_3^-)\}$ . When the network is trained we want the network’s output to  
 309 closely approximate the matrix  $R$ , that is we want  $\mathbf{y} \approx R\mathbf{x}$ . To do so, we first generate training  
 310 data of the form  $\{(\mathbf{x}, R\mathbf{x})\}$  by randomly sampling  $\mathbf{x}$  from the surface of a unit sphere, and train the  
 311 network using the procedure described in the previous section. To shorten the training time, we want  
 312 the magnitude of output  $\mathbf{y}$  to be of the same order as that of the input, therefore we make sure that the  
 313 maximum eigenvalue of  $R$  is close to one. We do this by first generating an arbitrary matrix  $R'$  with  
 314 random entries between -1 and 1, and then normalizing it by dividing it with maximum eigenvalue  
 315 :  $R = R' / \max\{eig(R')\}$ . Input and output data is generated using this matrix  $R$ . The network  
 316 is trained using this ideal data, meaning each training step sees an entirely new data point. In the  
 317 computer science community, this type of task is known as linear regression.

### 318 Appendix C: Analyzing the time dependence of the relaxation

319 To more rigorously ascertain whether the error curves depicted in Fig. 1 follow a non-exponential  
 320 relaxation, we analyzed their log-linear plots. As evident in Fig. A1, these plots deviate significantly  
 322 from a straight line, suggesting a departure from a simple exponential relationship. To demonstrate  
 323 the significance of the fit parameters, we examined their behavior with progressive increments in  
 324 the data size ( Fig. A2 ), fitting them to the non-exponential function. We observed that increasing  
 325 the data size leads the fit parameters to converge to a stable value, especially when the error curve  
 326 reaches saturation. In some instances, the error curve takes an extended period to saturate, resulting  
 327 in less stable fit parameters. This phenomenon is particularly noticeable for smaller  $TP$  values,  
 328 where the relaxation process is sluggish. Even with prolonged computations, the error curve does  
 329 not achieve saturation under these conditions.

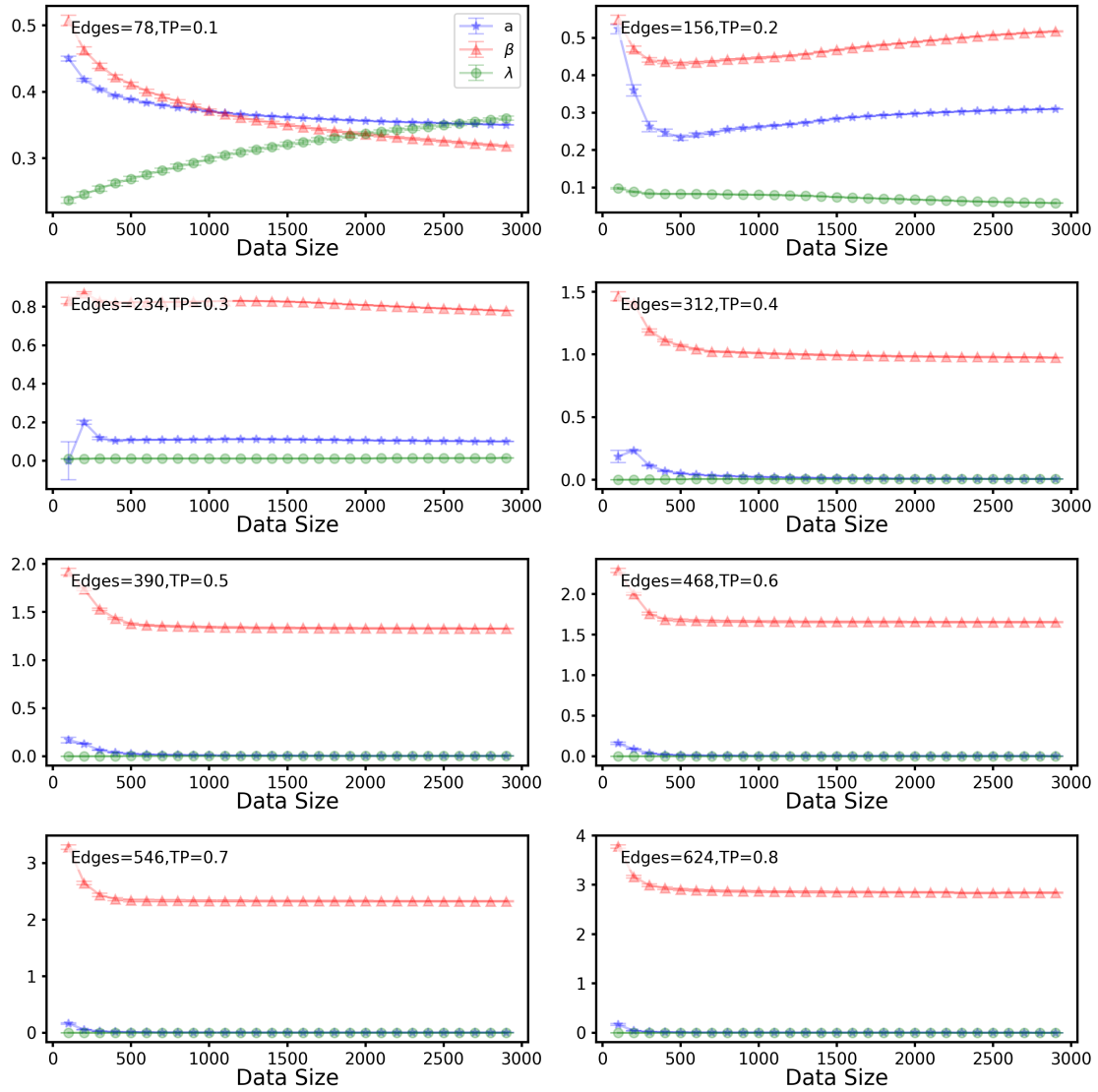


FIG. A2. *Fit Robustness Plots*: This figure illustrates the robustness of fit parameters depicted in Fig. 2. Each subplot represents the evolution of fit parameters with increasing data size. Notably, for  $TP$  values of 0.1 and 0.2, the fit parameters exhibit less stability. In contrast, higher  $TP$  values yield more consistent and stable parameters.

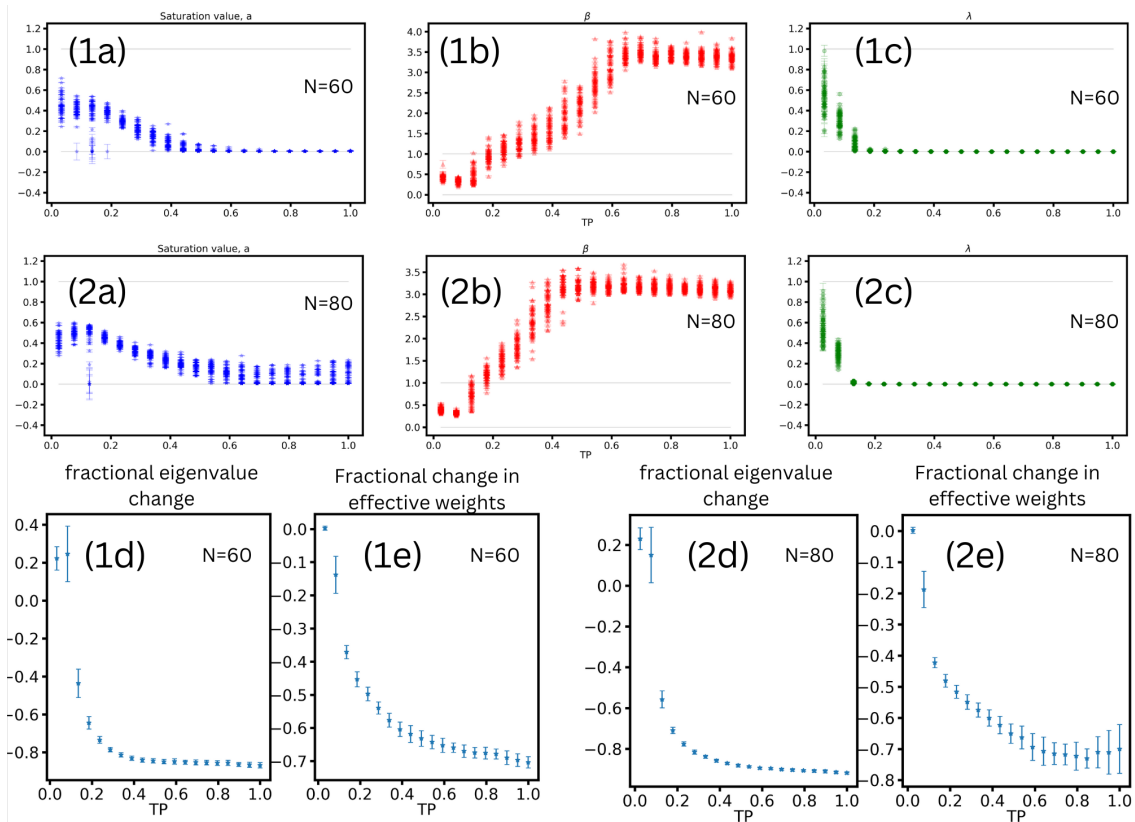


FIG. A3. Analysis for larger networks: [1a-c] shows fit parameters with increasing  $TP$  and [1d-e] shows fractional eigenvalue change and fractional change in effective weights for 60 node networks [2a-e]. Shows similar analysis for 80 node networks.

### 331 Appendix D: Analysis for larger network sizes

332 To further investigate the learning performance, we expanded our analysis from the previously studied  
 333 40-node network (as depicted in Fig. 2 and Fig. 3) to networks consisting of 60 and 80 nodes,  
 334 represented in Fig. A3. These networks were tasked with learning linear transformations of sizes  
 335  $15 \times 15$  and  $20 \times 20$ , respectively, while varying the number of edges. We observe a sharper  
 336 in  $\beta$  value for larger network sizes, but the saturation value of the error curve starts saturating above  
 337 zero.

### 339 Appendix E: The two-point correlation function for resistance networks

340 The effective resistance between nodes  $i$  and  $j$  in a resistance network, denoted as  $r_{ij}$ , provides a  
 341 measure of the ‘distance’ between these nodes [29]. It quantifies the potential difference between  
 342 nodes  $i$  and  $j$  when a unit current is injected at  $i$  and extracted at  $j$ , normalized by the total current.  
 343 The mathematical representation is:

$$r_{ij} = (\delta_i - \delta_j)^T L^+ (\delta_i - \delta_j) \quad (\text{A1})$$

344 In this expression,  $\delta_i$  and  $\delta_j$  are the Kronecker delta vectors. For a network with  $N$  nodes,  $\delta_i$  is a  
 345 vector of length  $N$  that has a value of 1 at the  $i$ -th position and 0 everywhere else. Similarly,  $\delta_j$   
 346 has a value of 1 at the  $j$ -th position and 0 elsewhere. . The term  $L^+$  represents the Moore-Penrose  
 347 pseudoinverse of the Laplacian matrix  $L$  of the network.

348 As we consider a network with a constant node count and progressively introduce more edges, the  
 349 average ‘distance’ between nodes diminishes. This is reflected in the reduction of the average path  
 350 length, which represents the mean number of steps required to traverse from one node to another,

351 averaged over all node pairs and paths. In the extreme case of a fully connected network, all nodes  
352 are adjacent. The effective resistance metric adeptly captures this behavior: as more edges are added,  
353 the effective resistance between nodes decreases.

354 Building on this understanding, we can conceptualize a two-point correlation function as:

$$p(r) = \frac{1}{N} \sum_{i \neq j} \delta(r - r_{ij}) \quad (\text{A2})$$

355 Where  $p(r)$  is normalized and  $N$  is chosen such that the area under  $p(r)$  is 1. In Fig. 4 we plot this  
356  $\langle p(r) \rangle$  averaged over 50 different network initializations.