

HOW CONNECTIVITY STRUCTURE SHAPES RICH AND LAZY LEARNING IN NEURAL CIRCUITS

Yuhan Helena Liu^{1,2,3,*}, Aristide Baratin⁴, Jonathan Cornford^{3,5}, Stefan Mihalas^{1,2}, Eric Shea-Brown^{1,2}, and Guillaume Lajoie^{3,6,7,*}

¹University of Washington, Seattle, WA, USA

²Allen Institute for Brain Science, Seattle WA, USA

³Mila - Quebec AI Institute, Montreal, QC, Canada

⁴Samsung - SAIT AI Lab, Montreal, QC, Canada

⁵McGill University, Montreal, QC, Canada

⁶Canada CIFAR AI Chair, CIFAR, Toronto, ON, Canada

⁷Université de Montréal, Montreal, QC, Canada

*Correspondence: hylu24@uw.edu, g.lajoie@umontreal.ca

ABSTRACT

In theoretical neuroscience, recent work leverages deep learning tools to explore how some network attributes critically influence its learning dynamics. Notably, initial weight distributions with small (resp. large) variance may yield a rich (resp. lazy) regime, where significant (resp. minor) changes to network states and representation are observed over the course of learning. However, in biology, neural circuit connectivity could exhibit a low-rank structure and therefore differs markedly from the random initializations generally used for these studies. As such, here we investigate how the structure of the initial weights — in particular their effective rank — influences the network learning regime. Through both empirical and theoretical analyses, we discover that high-rank initializations typically yield smaller network changes indicative of lazier learning, a finding we also confirm with experimentally-driven initial connectivity in recurrent neural networks. Conversely, low-rank initialization biases learning towards richer learning. Importantly, however, as an exception to this rule, we find lazier learning can still occur with a low-rank initialization that aligns with task and data statistics. Our research highlights the pivotal role of initial weight structures in shaping learning regimes, with implications for metabolic costs of plasticity and risks of catastrophic forgetting.

1 INTRODUCTION

Structural variations can significantly impact learning dynamics in theoretical neuroscience studies of animals. For instance, studies have revealed that specific neural connectivity patterns can facilitate faster learning of certain tasks (Braun et al., 2022; Raman & O’Leary, 2021; Simard et al., 2005; Canatar et al., 2021; Xie et al., 2022; Goudar et al., 2023; Chang et al., 2023). In deep learning, structure, encompassing architecture and initial connectivity, crucially dictates learning speed and effectiveness (Richards et al., 2019; Zador, 2019; Yang & Molano-Mazón, 2021; Braun et al., 2022).

A key structural aspect is the initial connectivity prior to training. Specifically, the initial connection weight magnitude can significantly bias learning dynamics, pushing them towards either rich or lazy regimes (Chizat et al., 2019; Flesch et al., 2021). Lazy learning often induces minor changes in the network during the learning process. Such minimal adjustments are advantageous given that plasticity is metabolically costly (Mery & Kawecky, 2005; Plaçais & Preat, 2013), and significant changes in representations might lead to issues like catastrophic forgetting (McCloskey & Cohen, 1989; Kirkpatrick et al., 2017). On the other hand, the rich learning regime can significantly adapt the network’s internal representations to task statistics, which can be advantageous for task feature acquisition and has implications for generalization (Flesch et al., 2021; George et al., 2022). Most research on initial weight magnitude’s role in learning dynamics has focused on random Gaussian

or Uniform initializations (Woodworth et al., 2020; Flesch et al., 2021; Braun et al., 2022). These patterns stand in contrast to the connectivity structures observed in biological neural circuits, which could exhibit a more pronounced low-rank eigenstructure (Song et al., 2005). This divergence prompts a pivotal question: how does the initial weight structure, given a fixed initial weight magnitude, bias the learning regime?

This study examines how initial weight structure, particularly the effective rank, modulates the *effective* richness or laziness of task learning within the standard training regime. We note that *rich* and *lazy* learning regimes have well established meanings in deep learning theory. The latter being defined as a situation where the Neural Tangent Kernel (NTK) stays stationary during training, while the former refers to the case where the NTK changes. In this work, we slightly extend these definitions and introduce **effective learning richness/laziness**. Unlike the traditional definition, which is based upon initialization, effective learning richness/laziness is defined in terms of post-training adjustment measurements. From this perspective, a learning process is deemed effectively "lazy" if the measured NTK movement is small. For example, consider a network whose initialization puts it in standard rich regime, but for a given task, its NTK moves very little during training. We define learning for this specific situation as effectively lazy. In other words, while the standard regime definition informs us (prior to training) whether the network can adapt significantly to task training or not, our "effective" definition lies in the post-training effects.

1.1 CONTRIBUTIONS

Our main **contributions** and findings can be summarized as follows:

- Through theoretical derivation in two-layer feedforward linear network, we demonstrate that higher-rank initialization results in *effectively* lazier learning **on average** across tasks (Theorem 1). We note that the emphasis of the theorem is on the expectation across tasks.
- We validate our theoretical findings in recurrent neural networks (RNNs) through numerical experiments on well-known neuroscience tasks (Figure 1) and demonstrate the applicability to different initial connectivity structures extracted from neuroscience data (Figure 2).
- We identify scenarios where certain low-rank initial weights still result in *effectively* lazier learning for specific tasks (Proposition 1 and Figure 3). We postulate that such patterns emerge when a neural circuit is predisposed — perhaps due to evolutionary factors or post-development — to certain tasks, ingraining specific inductive biases in neural circuits.

1.2 RELATED WORKS

An extended discussion on related works can also be found in Appendix A.

Theoretical Foundations of Neural Network Regimes and Implications for Neural Circuits:

The deep learning community has made tremendous strides in developing theoretical groundings for artificial neural networks (Advani et al., 2020; Jacot et al., 2018; Alemohammad et al., 2020; Agarwala et al., 2022; Atanasov et al., 2021; Azulay et al., 2021; Emami et al., 2021). A focal point is the 'rich' and 'lazy' learning regimes dichotomy, which have distinct impacts on representation and generalization (Chizat et al., 2019; Flesch et al., 2021; Geiger et al., 2020; George et al., 2022; Ghorbani et al., 2020; Woodworth et al., 2020; Paccolat et al., 2021; Nacson et al., 2022; HaoChen et al., 2021; Flesch et al., 2023). The 'lazy' regime results in minimal weight changes, while the 'rich' regime fosters task-specific adaptations. The transition between these is influenced by various factors, including initial weight scale and network width (Chizat et al., 2019; Geiger et al., 2020).

Deep learning theories increasingly inform studies of biological neural network learning dynamics (Bordelon & Pehlevan, 2022; Liu et al., 2022a; Braun et al., 2022; Ghosh et al., 2023; Saxe et al., 2019; Farrell et al., 2022; Pappan et al., 2020; Tishby & Zaslavsky, 2015). For the rich/lazy regime theory, the existence of diverse learning regimes in neural systems is evident through the resource-intensive plasticity-driven transformations prevalent in developmental phases, followed by more subdued adjustments (Lohmann & Kessels, 2014), and previous investigations characterized neural network behaviors under distinct regimes (Bordelon & Pehlevan, 2022; Schuessler et al., 2023) and discerning which mode yields solutions mimicking neural data (Flesch et al., 2021). Our work extends these studies by examining how initial weight structures affect learning.

Neural circuit initialization, connectivity patterns and learning: Extensive research has explored the influence of various random initializations on deep network learning (Saxe et al., 2013; Bahri et al., 2020; Glorot & Bengio, 2010; He et al., 2015; Arora et al., 2019). The literature predominantly focuses on random initialization, but actual neural structures exhibit markedly different connectivity patterns, such as Dale’s law and enriched cell-type-specific connectivity motifs (Rajan & Abbott, 2006; Ipsen & Peterson, 2020; Harris et al., 2022; Dahmen et al., 2020; Aljadeff et al., 2015). Motivated by existing evidence of low-rankedness in the brain (Thibeault et al., 2024) and the overrepresentation of local motifs in neural circuits (Song et al., 2005), which could be indicative of low-rank structures due to their influence on the eigenspectrum (Dahmen et al., 2020; Shao & Ostojic, 2023), our study explores the impact of connectivity effective rank on learning regimes. This focus is driven by the plausible presence of such low-rank structures in the brain, potentially revealed through these local motifs. With emerging connectivity data (Campagnola et al., 2022; MICrONS Consortium et al., 2021; Dorkenwald et al., 2022; Winnubst et al., 2019; Scheffer et al., 2020), future work is poised to encompass rich additional features of connectivity.

2 SETUP AND THEORETICAL FINDINGS

2.1 RNN SETUP

We examine recurrent neural networks (RNNs) because they are commonly adopted for modeling neural circuits (Barak, 2017; Song et al., 2016). We consider a RNN with N_{in} input units, N hidden units and N_{out} readout units (Figure 1A). The update formula for $h_t \in \mathbb{R}^N$ (the hidden state at time t) is governed by (Ehrlich et al., 2021; Molano-Mazon et al., 2022):

$$h_{t+1} = \rho h_t + (1 - \rho)(W_h f(h_t) + W_x x_t), \quad (1)$$

where an exponential Euler approximation is made with $\rho = e^{-dt/\tau_m} \in \mathbb{R}$ denoting the leak factor for simulation time step dt and τ_m denoting the membrane time constant; $f(\cdot) : \mathbb{R}^N \rightarrow \mathbb{R}^N$ is the activation function, for which we use *ReLU*; $W_h \in \mathbb{R}^{N \times N}$ (resp. $W_x \in \mathbb{R}^{N \times N_{in}}$) is the recurrent (resp. input) weight matrix and $x_t \in \mathbb{R}^{N_{in}}$ is the input at time step t . Readout $\hat{y}_t \in \mathbb{R}^{N_{out}}$, with readout weights $w \in \mathbb{R}^{N_{out} \times N}$, is defined as

$$\hat{y}_t = \langle w, f(h_t) \rangle. \quad (2)$$

The objective is to minimize scalar loss $L \in \mathbb{R}$, for which we use the cross-entropy loss for classification tasks and mean squared error for regression tasks. L is minimized by updating the parameters using variants of gradient descent:

$$\Delta W = -\eta \nabla_W L, \quad (3)$$

for learning rate $\eta \in \mathbb{R}$ and $W = [W_h \quad W_x \quad w^T] \in \mathbb{R}^{N \times (N_{in} + N + N_{out})}$ contains all the trainable parameters. Details of parameter settings can be found in Appendix C.

2.2 EFFECTIVE LAZINESS MEASURES

As mentioned above, we introduce *effective* richness and laziness, with effectively lazier (resp. richer) learning corresponding to less (resp. greater) network change over the course of learning. To quantify network change, we adopt the following three measures that have been used previously (George et al., 2022). We note that these measures can be sensitive to other architectural aspects that bias learning regimes, such as network width, so throughout we hold these variables constant when making the comparisons.

Weight change norm quantifies the vector norm of change in W . Effectively lazier learning should result in a lower weight change norm, and it is quantified as:

$$\|\Delta W\| := \|W^{(f)} - W^{(0)}\|, \quad (4)$$

where $\|\cdot\| = \|\cdot\|_F$; $W^{(0)}$ (resp. $W^{(f)}$) are the weights before (resp. after) training.

Representation alignment (RA) quantifies the directional change in a representational similarity matrix (RSM) before and after training. RSM focuses on the similarity between how two pairs

of input are represented by computing the Gram matrix R of last step hidden activity. Greater representation alignment indicates effectively lazier learning in the network, and it is obtained by

$$RA(R^{(f)}, R^{(0)}) := \frac{\text{Tr}(R^{(f)}R^{(0)})}{\|R^{(f)}\|\|R^{(0)}\|}, \quad \text{where } R := H^T H, \quad (5)$$

where $H \in \mathbb{R}^{N \times m}$ is the hidden activity at the last time step; $R^{(0)}$ and $R^{(f)} \in \mathbb{R}^{m \times m}$ are the initial and final RSM, respectively; m is the batch size.

Tangent kernel alignment (KA) quantifies the directional change in the neural tangent kernel (NTK) before and after training; effectively lazier learning should result in higher tangent kernel alignment. The NTK computes the Gram matrix K of the output gradient. Greater tangent kernel alignment points to effectively lazier learning, and it is obtained by

$$KA(K^{(f)}, K^{(0)}) := \frac{\text{Tr}(K^{(f)}K^{(0)})}{\|K^{(f)}\|\|K^{(0)}\|}, \quad \text{where } K := \nabla_W \hat{y}^T \nabla_W \hat{y} \quad (6)$$

where $K^{(0)}$ and $K^{(f)} \in \mathbb{R}^{m \times m}$ (for the $N_{out} = 1$ case) denote the initial and final NTK, respectively.

2.3 THEORETICAL FINDINGS

This subsection derives the theoretical impact of initial weight effective rank on tangent kernel alignment. First, Theorem 1 focuses on **task-agnostic** settings, treating task definition as random variables and computing the **expected** tangent kernel alignment across tasks. With some assumptions, tangent kernel alignment is maximized when the initial weight singular values are distributed across all dimensions (i.e. high-rank initialization).

In this section, our theoretical results are framed in a simplified feedforward setting, as we use a two-layer network with linear activations. However, we return to RNNs (Eq. 1) for the rest of the paper, and verify the generality of our theoretical findings with numerical experiments for both feedforward and recurrent architectures. Our choice is motivated by the need for theoretical tractability. While research on RNN learning in the NTK regime exists (Yang, 2020; Alemohammad et al., 2020; Emami et al., 2021), we are not aware of any studies featuring the final converged NTK that could serve as a basis for our comparison of the initial and final kernel. Consequently, we have chosen to focus on RNNs for neural circuit modeling and employ linear feedforward networks for theoretical derivations, a strategy also adopted by Farrell et al. (2022); numerous other studies, including Saxe et al. (2019), (Atanasov et al., 2021), (Arora et al., 2019), and (Braun et al., 2022), have similarly concentrated on extracting theoretical insights from linear feedforward networks.

For a two-layer linear network with input data $X \in \mathbb{R}^{d \times m}$, $W_1 \in \mathbb{R}^{N \times d}$ and $W_2 \in \mathbb{R}^{1 \times N}$ as weights for layers 1 and 2, respectively, the NTK throughout training, K , is:

$$K = X^T (W_1^T W_1 + \|W_2\|^2 I) X. \quad (7)$$

Without the loss of generality, suppose the output target $Y \in \mathbb{R}^{1 \times m}$ is generated from a linear teacher network as $Y = \beta^T X$, for some Gaussian vector $\beta \in \mathbb{R}^d$, with $\beta_i \sim \mathcal{N}(0, 1/d)$.

Theorem 1. (Informal) Consider the network above with its corresponding NTK in Eq. 7, trained under MSE loss with small initialization and whitened data. The expected kernel alignment across tasks is maximized with high-rank initialization, i.e. the singular values of $W_1^{(0)}$ are distributed across all dimensions. (Formal statement and proof are in Appendix B)

The intuition of Theorem 1 result is that, when two random vectors are drawn in high-dimensional spaces, corresponding to the low-rank initial network and the task, the probability of them being nearly orthogonal is very high; this then necessitates greater movement to eventually learn the task direction. We emphasize again that Theorem 1 is **task-agnostic**, i.e. it focuses on the **expected** tangent kernel alignment across input-output definitions. This is in contrast to **task-specific** settings (e.g. Woodworth et al. (2020)) that focus on a given task. In such task-specific settings, certain low-rank initializations can in fact lead to lazier learning. The following proposition predicts that if the task structure is known, low-rank initialization that is already aligned with the task statistics (input/output covariance) can lead to kernel alignment. We revisit this proposition again in Figure 3. We remark that initializing this way can still have high initial error because of randomized $W_2^{(0)}$.

Proposition 1. (Informal) Following the setup and assumptions in Theorem 1, rank-1 initializations of the form $W_1^{(0)} = \sigma[\beta^T / \|\beta\| \vec{0} \dots \vec{0}]$ leads to a high tangent kernel alignment. (Formal statement and proof are in Appendix B)

Above, we state technical results in terms of one metric of the effective laziness of learning — based on the NTK; our proof in Appendix B easily extends also to the representation alignment metric. The impact on weight change is also assessed in Appendix Proposition 2. This is in line with our simulations with RNNs, which will show similar trends for all three of the metrics introduced in Section 2.2).

3 SIMULATION RESULTS

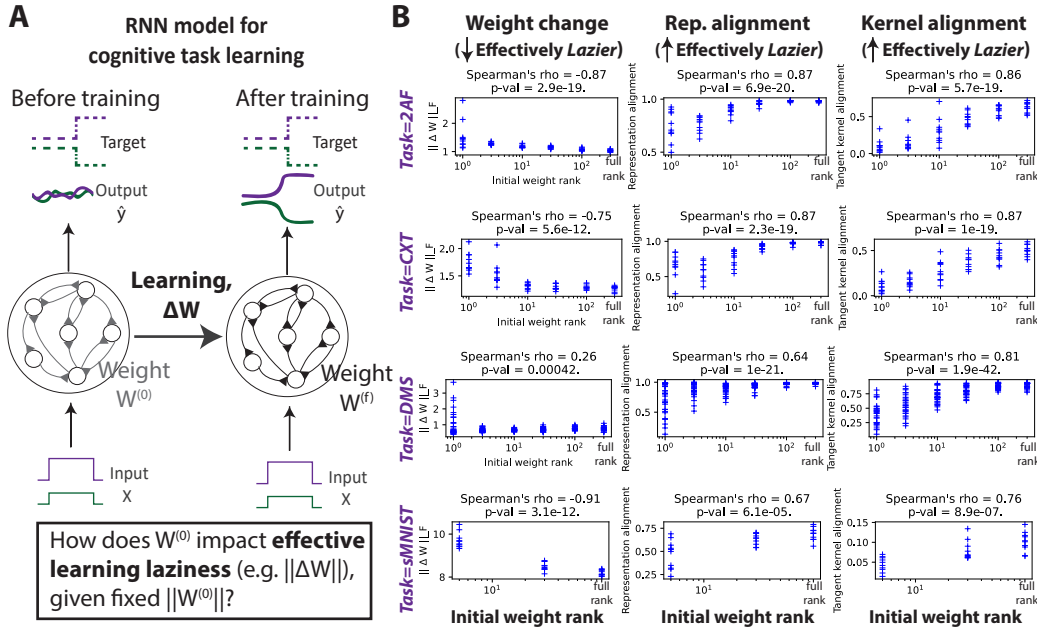


Figure 1: **Low-rank initial recurrent weights, generated using SVD, lead to greater changes (or effectively richer learning) in the recurrent neural network.** A) Schematic of RNN training setup. B) Measurements of effective richness vs laziness of learning (metrics as defined in Section 2.2), for RNN trained on several cognitive tasks in Neurogym (Molano-Mazon et al., 2022) as well as the sequential MNIST task (sMNIST). For details on SVD weight creation, see Appendix C. Fewer rank points were used for sMNIST due to computational time. Each dot represents a single training run, with each run using a different random initialization (10 runs total for each setting).

In this section we empirically illustrate and verify our main theoretical results, which are: (1) on average, high-rank initialization leads to effectively lazier learning (Theorem 1); (2) it is still possible for certain low-rank initializations that are already aligned to the task statistics to achieve effectively lazier learning (Proposition 1).

Impact on effective laziness by low-rank initialization via SVD in RNNs: As a proof-of-concept, we start in Figure 1 with low-rank initialization in RNNs by truncating an initial Gaussian random matrix via Singular Value Decomposition (SVD), which enables us to precisely control the rank, and rescale it to ensure that the comparison is across the same weight magnitude (Schuessler et al., 2020). Additionally, all comparisons were made after training was completed, and all these training sessions achieved comparable losses. For our investigations, we applied this initialization scheme across a variety of cognitive tasks — including two-alternative forced choice (2AF), delayed-match-to-sample (DMS), context-dependent decision-making (CXT) tasks — implemented with Neurogym (Molano-Mazon et al., 2022) and the well-known machine learning benchmark sequential MNIST (sMNIST). Figure 1 indicates that low-rank initial weights result in effectively richer learning and greater network changes.

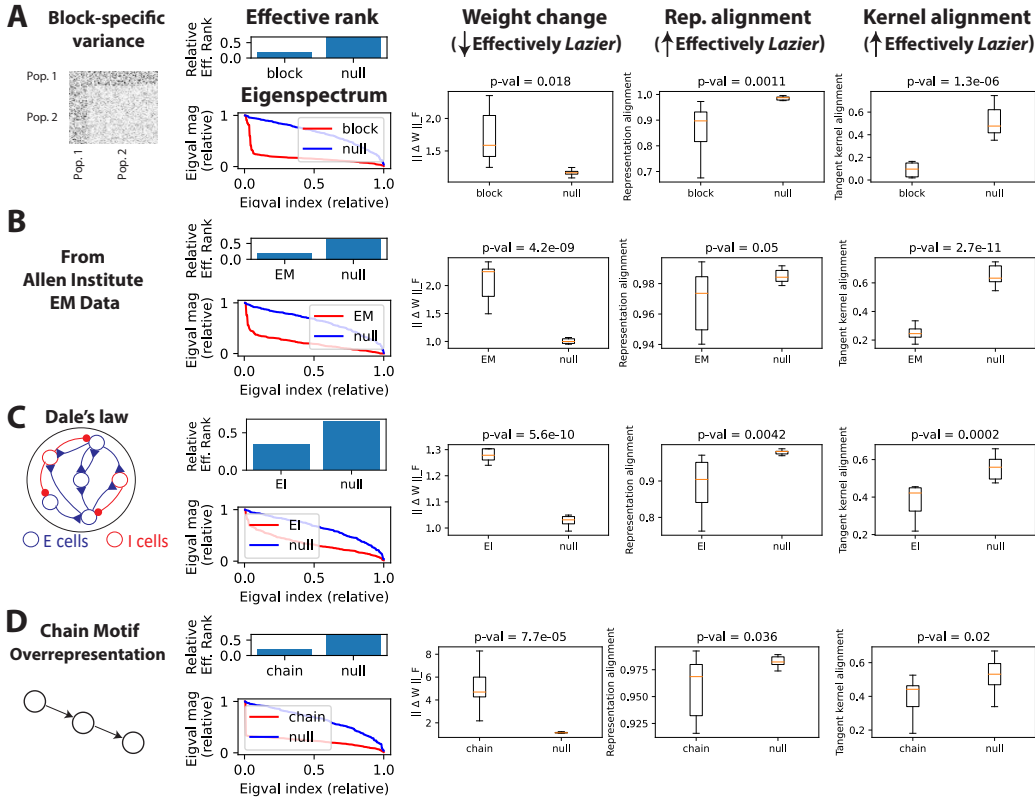


Figure 2: Low-rank initial weight structures, inspired by biological examples, lead to effectively richer learning. We present the eigenspectrum and the relative effective rank of connectivity in A) structures with cell-type-specific statistics, B) structures derived from EM data, C) structures obeying Dale’s law, and D) structures with an over-representation of chain motifs; we also present the effective learning laziness for networks initialized with these connectivity structures. These structures exhibit a lower effective rank compared to standard random Gaussian initialization (null). We plotted the magnitude of the eigenvalues (Eigval mag) — scaled by the dominant eigenvalue’s magnitude — against their indices normalized by the network size N (Eigval index). We apply the effective laziness measures described in Section 2.2 to compare the effective laziness of experimentally-driven initial connectivity versus standard random Gaussian initialization (null). See Appendix C for details on network initialization. The boxplots are generated from 10 independent runs with different initialization seeds. Due to space constraints, we include only the 2AF task here, but Appendix Figures 5 and 6 show that similar trends hold for the DMS and CXT tasks.

These numerical trends are in line with Theorem 1, which focused on an idealized setting of a two-layer linear network with numerical results in Appendix Figure 4A. We also demonstrated this trend for a non-idealized feedforward setting in Appendix Figure 4B, and more explorations in feedforward settings and across a broader range of architecture is left for future exploration due to our focus on RNNs. In the Appendix, we show the main trends observed in Figure 1 also hold for Uniform initialization (Figure 7), soft initial weight rank (Figure 8), various network sizes (Figure 9), learning rates (Figure 10), gains (Figure 11), and finer time step dt (Figure 12). We note that, in addition to fixing the weight magnitude across comparisons, the dynamical regime might also confound learning regimes. A common method for controlling the dynamical regime is through the leading weight eigenvalue, which affects the top Lyapunov exponent. Controlling in this manner led to similar trends (Appendix Figure 13). Investigating the relationship between learning regimes and various concepts of dynamical regimes further is a promising direction for future work. Moreover, since our emphasis is on the effective learning regime, which is based on post-training changes, we concentrated on the laziness measures computed from networks after training, rather than during the learning process. However, we also tracked the alignment with the initial kernel and task kernel alignment during

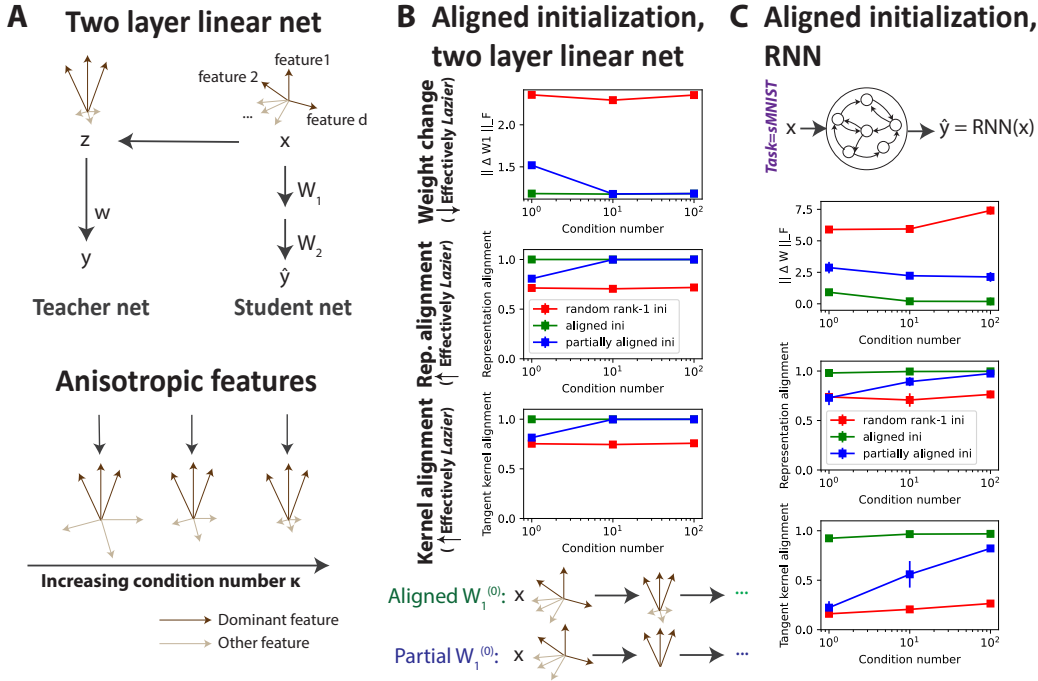


Figure 3: **Low-rank initializations can still achieve high alignment for specific tasks (see Proposition 1).** A) The student-teacher two-layer linear network setup as described in Section 2.3, but with feature anisotropy controlled by a feature modulation matrix F , i.e. $z = Fx$. The condition number of F dictates the relative feature strength. We set the top half of the singular values of F are set to κ , while the bottom half are set to 1, where κ represents the condition number of F . B) The aligned initialization (green) is achieved by setting W_1 as described in Proposition 1 (with $\beta = w^T F$, w is as illustrated), so that the initialization aligns with the task statistics. The partial alignment (blue) mirrors the aligned case, but F is substituted with its rank- $(d/2)$ truncation, causing the network to align only with the dominant features. **We observe that a considerably higher alignment can be achieved when the initialization aligns solely with the dominant features, especially when the relative strength of these dominant features is high.** C) The analysis from B) is replicated for RNNs learning the sMNIST task. As the ground truth network function is elusive, we use a teacher network with pre-trained weights. Once again, we replace F with its rank- $(d/2)$ truncation for partial alignment. Details on the input/output definitions and initializations, as well as other simulation specifics, are available in Appendix C. We note that in all scenarios presented here, the initial errors are high since the readout weights are initialized randomly, rendering it a valid learning problem.

training (Appendix Figure 14). We also examined how the kernel’s effective rank evolves throughout the training period (Appendix Figure 15).

Low-rank initialization via biologically motivated connectivity in RNNs: To establish a closer connection with biological neural circuits, we have tested our predictions on low-rank initialization using a variety of biologically motivated structures capable of resulting in low-rank connectivity. Here are some of the examples: (A) connectivity with cell-type-specific statistics (Aljadeff et al., 2015), where each block in the weight matrix corresponds to the connections between neurons of two distinct cell types, with the variance of these connections differing from one block to another. In terms of block-specific connectivity statistics, there are infinite possibilities for defining the blocks, each resulting in a unique eigenspectrum. For the example provided here, we adopted the setup from Figure S3 in Aljadeff et al. (2015), with parameters set as $\alpha = 0.02$, $\gamma = 10$, and $1 - \epsilon = 0.8$; these correspond to the fraction of hyperexcitable neurons, gain of hyperexcitable connections and gain of the rest, respectively. We follow this particular setup because it has been demonstrated to create an outlier leading eigenvalue, thereby reducing the effective rank. We also consider (B) connectivity matrix derived from the electron microscopy (EM) data (Allen Institute, 2023), where the synaptic connections between individual neurons are meticulously mapped to create a detailed

and comprehensive representation of neural circuits. Also, we consider (C) connectivity obeying Dale’s law, where each neuron is either excitatory or inhibitory, meaning it can only send out one type of signal – either increasing or decreasing the activity of connected neurons – a principle inspired by the way neurons behave in biological systems (Song et al., 2005). Additionally, (D) the over-representation of certain localized connectivity patterns (or network motifs) — such as the chain motif, where two cells are connected via a third intermediary cell — creates outliers in the weight eigenspectrum, subsequently lowering the effective rank (Zhao et al., 2011; Hu et al., 2018; Dahmen et al., 2020). Details of these initial connectivity structures are provided in Appendix C.

As illustrated in Figure 2, these connectivity structures, motivated by known features of biological neural networks, exhibit a lower effective rank compared to standard random Gaussian initialization, thereby serving as natural testbeds for our theoretical predictions. To quantify (relative) effective rank, we used $(\sum_i |\lambda_i|)/(|\lambda_1|N)$, which indicates the fraction of eigenvalues on the order of the dominant one and captures the (scaled) area under the curve of the eigenspectrum plots. We also tried effective rank based on singular values, i.e. $(\sum_i |s_i|)/(|s_1|N)$, in Appendix Figure 16 and observed similar trends. Importantly, Figure 2 show that these different low-rank biologically motivated structures can lead to effectively richer learning compared to the standard random Gaussian initialization. This finding supports our overarching prediction, that lower rank initial weights leads to effectively richer learning. We note that to test our theoretical predictions based on gradient-descent learning without specific constraints on the solutions, the structures are enforced only at initialization and not constrained during training. In Appendix Figure 17, we also constrained Dale’s Law throughout training and found similar trends.

Low-rank initialization aligned with task statistics: These simulations may be considered to be within our task-agnostic framework. That is, we have chosen a “random” battery of tasks that is not directly matched to the initial network connectivity structures. Thus, our findings that lower rank initializations lead to richer learning are expected from our theoretical prediction on the task-averaged alignment (Theorem 1), rather than something task-specific. However, Proposition 1 also predicts that low-rank initialization can lead to lazy learning if the initialization is already aligned to the task structure. To test this, we observe in Figure 3 that a considerably higher alignment can be achieved when the initialization aligns solely with the dominant task features, especially when the relative strength of these dominant features is high. We postulate that such alignment may occur in biological settings if the circuit has evolved to preferentially learn specific tasks.

4 DISCUSSION

Our investigation casts light on the nuanced influence of initial weight effective rank on learning dynamics. Anchored by Theorem 1, our theoretical findings underscore that high-rank random initialization generally facilitates *effectively* lazier learning on average across tasks. This focus on the expectation across tasks can provide insights into the circuit’s flexibility in learning across a broad range of tasks as well as predict the effective learning regime when the task structure is uncertain. However, certain low-rank initial weights, when naturally predisposed to specific tasks, may lead to effectively lazier learning, suggesting an interesting interplay between evolutionary or developmental biases and learning dynamics (Proposition 1). Our numerical experiments on RNNs further validate these theoretical findings illustrating the impact of initial rank in diverse settings.

Potential implications to neuroscience: We investigate the impact of effective weight rank on learning regimes due to its relevance in neuroscience. Learning regimes reflect the extent of change through learning, implicating metabolic costs and catastrophic forgetting (McCloskey & Cohen, 1989; Plaçais & Preat, 2013; Mery & Kawecki, 2005). The presence of different learning regimes is demonstrated in neural systems, since during developmental phases where neural circuits undergo extensive, plasticity-driven transformations. In contrast, mature neural circuits exhibit more subtle synaptic adjustments (Lohmann & Kessels, 2014). We hypothesize that a circuit’s task-specific alignment might be established either evolutionarily or during early development. The specialization of neural circuits, such as ventral versus dorsal (Bakhtiari et al., 2021), may arise from engaging in tasks with similar computational demands. Conversely, circuits with high-rank structures may be less specialized, handling a wider array of tasks. Our framework could be used to compare connectivities across brain regions and species in order to predict their function and flexibility, assessing their functional specialization based on effective connectivity rank. Additionally, our framework predicts

that connectivity rank will affect the degree of change in neural activity during the learning of new tasks. This hypothesis could be tested through BCI experiments, as shown in Sadtler et al. (2014) and Golub et al. (2018), to explore how learning dynamics vary with connectivity rank.

Regarding deep learning, low-rank initialization is not a common practice, yet adaptations like low-rank updates have gained popularity in training large models (Hu et al., 2021). LoRA, the study cited, concentrates on parameter updates rather than initializations, but understanding how update rank affects learning regimes is crucial. Our results offer a starting point for further exploration in this area. Although different rank initializations are less explored, with some exceptions like Vodrahalli et al. (2022), our findings suggest that this area should receive more attention due to its potential effects on learning regimes and, consequently, on generalization (George et al., 2022).

Limitations and future directions: Our study predominantly focused on the weight (effective) rank, leaving the exploration of other facets of weight on the effective learning regime as an open avenue. Also, the ramifications of effective learning regimes on learning speed — given the known results on kernel alignment and ease of learning (Bartlett et al., 2021) and present mixed findings in the existing literature (Flesch et al., 2021; George et al., 2022) — warrant further exploration.

Expanding the scope of our study calls for examining a wider variety of tasks, neural network architectures, and learning rules. Although our work is based on the backpropagation learning rule, its implications for biologically plausible learning rules remain unexplored. Our primary criterion for selecting tasks was their relevance to neuroscience, aligning with our main objectives. However, given the diverse range of tasks performed by various species, future research could benefit from exploring a more extensive array of tasks. Exploring more complex neuroscience tasks, such as those in Mod-Cog (Khona et al., 2023), could provide valuable insights. On that note, we tested the pattern generation task from Bellec et al. (2020), a neuroscience task differing in structure from the Neurogym tasks, and observed similar trends (refer to Appendix Figure 18).

Additionally, we ensured the consistency of outcomes against factors like width, learning rate, and initial gain (see Appendix D), but other factors such as dynamical regime and noise (HaoChen et al., 2021) remain underexamined. On that note, the study’s focus on RNNs with finite task duration prompts further investigation into the implications for tasks with extended time steps and how conclusions for feedforward network depth (Xiao et al., 2020; Seleznova & Kutyniok, 2022) translate to RNN sequence length. Examining several mechanisms at once is beyond the scope of one paper, but our theoretical work constitutes the foundation for future investigations.

Moreover, it is crucial to further explore the neuroscientific implications of effective learning regimes, as well as their diverse impacts on aspects such as representation, including kernel-task alignment (see Appendix Figure 14), and generalization capabilities (Flesch et al., 2021; George et al., 2022; Schuessler et al., 2023). Our current study did not delve into how initial weight rank affects these facets of learning, representing an essential future direction in connecting weight rank to these theoretical implications in both biological and artificial neural networks.

Furthermore, while there exists evidence for low-rankedness in the brain (Thibeault et al., 2024), the extent to which the brain uses low-rank structures remains an open question, especially as neural circuit structures can vary across regions and species. While local connectivity statistics (Song et al., 2005) can offer some predictive insight into the global low-rank structure, this relationship is not always immediately apparent (Shao & Ostojic, 2023). Our theoretical results contribute to understanding the role of connectivity rank in the brain by linking effective connectivity rank with learning dynamics.

Lastly, we have primarily examined low-rank tasks and there remains unexplored terrain regarding the interplay between the number of task classes and weight rank, which is pivotal to uncovering a more precise relationship between the effective learning regime and the initial weight rank (Dubreuil et al., 2022; Gao et al., 2017). Overall, this dynamic area of learning regimes is ripe for many explorations, integrating numerous factors; our work contributes to this exciting area with new tools.

5 ACKNOWLEDGEMENT

We thank Andrew Saxe, Stefano Recanatesi, Kyle Aitken and Dana Mastrovito for insightful discussions and helpful feedback. This research was supported by NSERC PGS-D (Y.H.L.); FRQNT B2X

(Y.H.L.); Pearson Fellowship (Y.H.L.); NSF AccelNet IN-BIC program, Grant No. OISE-2019976 AM02 (Y.H.L.); NIH BRAIN, Grant No. R01 1RF1DA055669 (Y.H.L., E.S.B., S.M.); Mitacs Globalink Research Award (Y.H.L.); IVADO Postdoctoral Fellowship (J.C); the Canada First Research Excellence Fund (J.C.); NSERC Discovery Grant RGPIN-2018-04821 (G.L); Canada Research Chair in Neural Computations and Interfacing (G.L.); Canada CIFAR AI Chair program (G.L.). We also thank the Allen Institute founder, Paul G. Allen, for his vision, encouragement, and support.

REFERENCES

- Madhu S Advani, Andrew M Saxe, and Haim Sompolinsky. High-dimensional dynamics of generalization error in neural networks. *Neural Networks*, 132:428–446, 2020.
- Atish Agarwala, Fabian Pedregosa, and Jeffrey Pennington. Second-order regression models exhibit progressive sharpening to the edge of stability. *arXiv preprint arXiv:2210.04860*, 2022.
- Sina Alemohammad, Zichao Wang, Randall Balestriero, and Richard Baraniuk. The recurrent neural tangent kernel. *arXiv preprint arXiv:2006.10246*, 2020.
- Johnatan Aljadeff, Merav Stern, and Tatyana Sharpee. Transition to chaos in random networks with cell-type-specific connectivity. *Physical review letters*, 114(8):088101, 2015.
- Allen Institute. Allen institute for brain science, 2023. <https://portal.brain-map.org/explore/connectivity/ultrastructural-connectomics/>.
- Raman Arora, Sanjeev Arora, Joan Bruna, Nadav Cohen, Simon Du, Rong Ge, Suriya Gunasekar, Chi Jin, Jason Lee, Tengyu Ma, et al. Theory of deep learning, 2019.
- Alexander Atanasov, Blake Bordelon, and Cengiz Pehlevan. Neural networks as kernel learners: The silent alignment effect. *arXiv preprint arXiv:2111.00034*, 2021.
- Shahar Azulay, Edward Moroshko, Mor Shpigel Nacson, Blake E Woodworth, Nathan Srebro, Amir Globerson, and Daniel Soudry. On the implicit bias of initialization shape: Beyond infinitesimal mirror descent. In *International Conference on Machine Learning*, pp. 468–477. PMLR, 2021.
- Yasaman Bahri, Jonathan Kadmon, Jeffrey Pennington, Sam S Schoenholz, Jascha Sohl-Dickstein, and Surya Ganguli. Statistical mechanics of deep learning. *Annual Review of Condensed Matter Physics*, 11:501–528, 2020.
- Shahab Bakhtiari, Patrick Mineault, Timothy Lillicrap, Christopher Pack, and Blake Richards. The functional specialization of visual cortex emerges from training parallel pathways with self-supervised predictive learning. *Advances in Neural Information Processing Systems*, 34: 25164–25178, 2021.
- Omri Barak. Recurrent neural networks as versatile tools of neuroscience research. *Current opinion in neurobiology*, 46:1–6, 2017.
- Aristide Baratin, Thomas George, César Laurent, R Devon Hjelm, Guillaume Lajoie, Pascal Vincent, and Simon Lacoste-Julien. Implicit regularization via neural feature alignment. In *International Conference on Artificial Intelligence and Statistics*, pp. 2269–2277. PMLR, 2021.
- Peter L Bartlett, Andrea Montanari, and Alexander Rakhlin. Deep learning: a statistical viewpoint. *Acta numerica*, 30:87–201, 2021.
- Guillaume Bellec, Franz Scherr, Anand Subramoney, Elias Hajek, Darjan Salaj, Robert Legenstein, and Wolfgang Maass. A solution to the learning dilemma for recurrent networks of spiking neurons. *Nature communications*, 11(1):3625, 2020.
- Blake Bordelon and Cengiz Pehlevan. The influence of learning rule on representation dynamics in wide neural networks. *arXiv preprint arXiv:2210.02157*, 2022.
- Lukas Braun, Clémentine Dominé, James Fitzgerald, and Andrew Saxe. Exact learning dynamics of deep linear networks with prior knowledge. *Advances in Neural Information Processing Systems*, 35:6615–6629, 2022.

- Luke Campagnola, Stephanie C Seeman, Thomas Chartrand, Lisa Kim, Alex Hoggarth, Clare Gamlin, Shinya Ito, Jessica Trinh, Pasha Davoudian, Cristina Radaelli, et al. Local connectivity and synaptic dynamics in mouse and human neocortex. *Science*, 375(6585):eabj5861, 2022.
- Abdulkadir Canatar, Blake Bordelon, and Cengiz Pehlevan. Spectral bias and task-model alignment explain generalization in kernel regression and infinitely wide neural networks. *Nature communications*, 12(1):2914, 2021.
- Joanna C Chang, Matthew G Perich, Lee E Miller, Juan A Gallego, and Claudia Clopath. De novo motor learning creates structure in neural activity space that shapes adaptation. *bioRxiv*, pp. 2023–05, 2023.
- Lenaic Chizat, Edouard Oyallon, and Francis Bach. On lazy training in differentiable programming. *Advances in neural information processing systems*, 32, 2019.
- David Dahmen, Stefano Recanatesi, Gabriel K Ocker, Xiaoxuan Jia, Moritz Helias, and Eric Shea-Brown. Strong coupling and local control of dimensionality across brain areas. *Biorxiv*, pp. 2020–11, 2020.
- Sigurd Diederich and Manfred Opper. Learning of correlated patterns in spin-glass networks by local learning rules. *Physical review letters*, 58(9):949, 1987.
- Sven Dorkenwald, Claire E McKellar, Thomas Macrina, Nico Kemnitz, Kisuk Lee, Ran Lu, Jingpeng Wu, Sergiy Popovych, Eric Mitchell, Barak Nehoran, et al. Flywire: online community for whole-brain connectomics. *Nature methods*, 19(1):119–128, 2022.
- Alexis Dubreuil, Adrian Valente, Manuel Beiran, Francesca Mastrogiuseppe, and Srdjan Ostojic. The role of population structure in computations through neural dynamics. *Nature neuroscience*, 25(6):783–794, 2022.
- Daniel B Ehrlich, Jasmine T Stone, David Brandfonbrener, Alexander Atanasov, and John D Murray. Psychrnn: An accessible and flexible python package for training recurrent neural network models on cognitive tasks. *Eneuro*, 8(1), 2021.
- Melikasadat Emami, Mojtaba Sahraee-Ardakan, Parthe Pandit, Sundeep Rangan, and Alyson K Fletcher. Implicit bias of linear rnns. In *International Conference on Machine Learning*, pp. 2982–2992. PMLR, 2021.
- Matthew Farrell, Stefano Recanatesi, Timothy Moore, Guillaume Lajoie, and Eric Shea-Brown. Gradient-based learning drives robust representations in recurrent neural networks by balancing compression and expansion. *Nature Machine Intelligence*, 4(6):564–573, 2022.
- Timo Flesch, Keno Juechems, Tsvetomira Dumbalska, Andrew Saxe, and Christopher Summerfield. Rich and lazy learning of task representations in brains and neural networks. *BioRxiv*, pp. 2021–04, 2021.
- Timo Flesch, Andrew Saxe, and Christopher Summerfield. Continual task learning in natural and artificial agents. *Trends in Neurosciences*, 46(3):199–210, 2023.
- Peiran Gao, Eric Trautmann, Byron Yu, Gopal Santhanam, Stephen Ryu, Krishna Shenoy, and Surya Ganguli. A theory of multineuronal dimensionality, dynamics and measurement. *BioRxiv*, pp. 214262, 2017.
- 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.
- Thomas George, Guillaume Lajoie, and Aristide Baratin. Lazy vs hasty: linearization in deep networks impacts learning schedule based on example difficulty. *TMLR*, 2022.
- Behrooz Ghorbani, Song Mei, Theodor Misiakiewicz, and Andrea Montanari. When do neural networks outperform kernel methods? *Advances in Neural Information Processing Systems*, 33:14820–14830, 2020.

- Arna Ghosh, Yuhan Helena Liu, Guillaume Lajoie, Konrad Kording, and Blake Aaron Richards. How gradient estimator variance and bias impact learning in neural networks. In *The Eleventh International Conference on Learning Representations*, 2023.
- 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*, pp. 249–256. JMLR Workshop and Conference Proceedings, 2010.
- Matthew D Golub, Patrick T Sadtler, Emily R Oby, Kristin M Quick, Stephen I Ryu, Elizabeth C Tyler-Kabara, Aaron P Batista, Steven M Chase, and Byron M Yu. Learning by neural reassociation. *Nature neuroscience*, 21(4):607–616, 2018.
- Vishwa Goudar, Barbara Peysakhovich, David J Freedman, Elizabeth A Buffalo, and Xiao-Jing Wang. Schema formation in a neural population subspace underlies learning-to-learn in flexible sensorimotor problem-solving. *Nature Neuroscience*, 26(5):879–890, 2023.
- Will Greedy, Heng Wei Zhu, Joseph Pemberton, Jack Mellor, and Rui Ponte Costa. Single-phase deep learning in cortico-cortical networks. *Advances in Neural Information Processing Systems*, 35:24213–24225, 2022.
- Jeff Z HaoChen, Colin Wei, Jason Lee, and Tengyu Ma. Shape matters: Understanding the implicit bias of the noise covariance. In *Conference on Learning Theory*, pp. 2315–2357. PMLR, 2021.
- Isabelle D Harris, Hamish Meffin, Anthony N Burkitt, and Andre DH Peterson. Eigenvalue spectral properties of sparse random matrices obeying dale’s law. *arXiv preprint arXiv:2212.01549*, 2022.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *Proceedings of the IEEE international conference on computer vision*, pp. 1026–1034, 2015.
- Geoffrey Hinton. The forward-forward algorithm: Some preliminary investigations. *arXiv preprint arXiv:2212.13345*, 2022.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021.
- Yu Hu, Steven L Brunton, Nicholas Cain, Stefan Mihalas, J Nathan Kutz, and Eric Shea-Brown. Feedback through graph motifs relates structure and function in complex networks. *Physical Review E*, 98(6):062312, 2018.
- Jesper R Ipsen and Andre DH Peterson. Consequences of dale’s law on the stability-complexity relationship of random neural networks. *Physical Review E*, 101(5):052412, 2020.
- Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization in neural networks. *Advances in neural information processing systems*, 31, 2018.
- Mikhail Khona, Sarthak Chandra, Joy J Ma, and Ila R Fiete. Winning the lottery with neural connectivity constraints: Faster learning across cognitive tasks with spatially constrained sparse rnns. *Neural Computation*, 35(11):1850–1869, 2023.
- James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. Overcoming catastrophic forgetting in neural networks. *Proceedings of the national academy of sciences*, 114(13):3521–3526, 2017.
- Axel Laborieux and Friedemann Zenke. Holomorphic equilibrium propagation computes exact gradients through finite size oscillations. *arXiv preprint arXiv:2209.00530*, 2022.
- Yann LeCun. The mnist database of handwritten digits. <http://yann.lecun.com/exdb/mnist/>, 1998.
- Timothy P Lillicrap, Adam Santoro, Luke Marris, Colin J Akerman, and Geoffrey Hinton. Backpropagation and the brain. *Nature Reviews Neuroscience*, 21(6):335–346, 2020.

- Yuhan Helena Liu, Stephen Smith, Stefan Mihalas, Eric Shea-Brown, and Uygur Sümbül. Cell-type-specific neuromodulation guides synaptic credit assignment in a spiking neural network. *Proceedings of the National Academy of Sciences*, 118(51):e2111821118, 2021.
- Yuhan Helena Liu, Arna Ghosh, Blake Richards, Eric Shea-Brown, and Guillaume Lajoie. Beyond accuracy: generalization properties of bio-plausible temporal credit assignment rules. *Advances in Neural Information Processing Systems*, 35:23077–23097, 2022a.
- Yuhan Helena Liu, Stephen Smith, Stefan Mihalas, Eric Shea-Brown, and Uygur Sümbül. Biologically-plausible backpropagation through arbitrary timespans via local neuromodulators. *arXiv preprint arXiv:2206.01338*, 2022b.
- Christian Lohmann and Helmut W Kessels. The developmental stages of synaptic plasticity. *The Journal of physiology*, 592(1):13–31, 2014.
- Owen Marschall, Kyunghyun Cho, and Cristina Savin. A unified framework of online learning algorithms for training recurrent neural networks. *The Journal of Machine Learning Research*, 21(1):5320–5353, 2020.
- Michael McCloskey and Neal J Cohen. Catastrophic interference in connectionist networks: The sequential learning problem. In *Psychology of learning and motivation*, volume 24, pp. 109–165. Elsevier, 1989.
- Frederic Mery and Tadeusz J Kawecki. A cost of long-term memory in drosophila. *Science*, 308(5725):1148–1148, 2005.
- Alexander Meulemans, Nicolas Zucchet, Seijin Kobayashi, Johannes Von Oswald, and João Sacramento. The least-control principle for local learning at equilibrium. *Advances in Neural Information Processing Systems*, 35:33603–33617, 2022.
- MICrONS Consortium, J Alexander Bae, Mahaly Baptiste, Caitlyn A Bishop, Agnes L Bodor, Derrick Brittain, JoAnn Buchanan, Daniel J Bumbarger, Manuel A Castro, Brendan Celi, et al. Functional connectomics spanning multiple areas of mouse visual cortex. *BioRxiv*, pp. 2021–07, 2021.
- Manuel Molano-Mazon, Joao Barbosa, Jordi Pastor-Ciurana, Marta Fradera, Ru-Yuan Zhang, Jeremy Forest, Jorge del Pozo Lerida, Li Ji-An, Christopher J Cueva, Jaime de la Rocha, et al. Neurogym: An open resource for developing and sharing neuroscience tasks. 2022.
- James M Murray. Local online learning in recurrent networks with random feedback. *Elife*, 8:e43299, 2019.
- Mor Shpigel Nacson, Kavya Ravichandran, Nathan Srebro, and Daniel Soudry. Implicit bias of the step size in linear diagonal neural networks. In *International Conference on Machine Learning*, pp. 16270–16295. PMLR, 2022.
- Jonas Paccolat, Leonardo Petrini, Mario Geiger, Kevin Tyloo, and Matthieu Wyart. Geometric compression of invariant manifolds in neural networks. *Journal of Statistical Mechanics: Theory and Experiment*, 2021(4):044001, 2021.
- Vardan Papyan, XY Han, and David L Donoho. Prevalence of neural collapse during the terminal phase of deep learning training. *Proceedings of the National Academy of Sciences*, 117(40):24652–24663, 2020.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32, 2019.
- Alexandre Payeur, Jordan Guerguiev, Friedemann Zenke, Blake A Richards, and Richard Naud. Burst-dependent synaptic plasticity can coordinate learning in hierarchical circuits. *Nature neuroscience*, pp. 1–10, 2021.

- Mohammad Pezeshki, Oumar Kaba, Yoshua Bengio, Aaron C Courville, Doina Precup, and Guillaume Lajoie. Gradient starvation: A learning proclivity in neural networks. *Advances in Neural Information Processing Systems*, 34, 2021.
- Pierre-Yves Plaçais and Thomas Preat. To favor survival under food shortage, the brain disables costly memory. *Science*, 339(6118):440–442, 2013.
- Roman Pogodin, Jonathan Cornford, Arna Ghosh, Gauthier Gidel, Guillaume Lajoie, and Blake Richards. Synaptic weight distributions depend on the geometry of plasticity. *arXiv preprint arXiv:2305.19394*, 2023.
- Kanaka Rajan and Larry F Abbott. Eigenvalue spectra of random matrices for neural networks. *Physical review letters*, 97(18):188104, 2006.
- Dhruva V Raman and Timothy O’Leary. Frozen algorithms: how the brain’s wiring facilitates learning. *Current Opinion in Neurobiology*, 67:207–214, 2021.
- Blake A Richards, Timothy P Lillicrap, Philippe Beaudoin, Yoshua Bengio, Rafal Bogacz, Amelia Christensen, Claudia Clopath, Rui Ponte Costa, Archy de Berker, Surya Ganguli, et al. A deep learning framework for neuroscience. *Nature neuroscience*, 22(11):1761–1770, 2019.
- Pieter R Roelfsema and Anthony Holtmaat. Control of synaptic plasticity in deep cortical networks. *Nature Reviews Neuroscience*, 19(3):166–180, 2018.
- João Sacramento, Rui Ponte Costa, Yoshua Bengio, and Walter Senn. Dendritic cortical microcircuits approximate the backpropagation algorithm. *arXiv preprint arXiv:1810.11393*, 2018.
- Patrick T Sadtler, Kristin M Quick, Matthew D Golub, Steven M Chase, Stephen I Ryu, Elizabeth C Tyler-Kabara, Byron M Yu, and Aaron P Batista. Neural constraints on learning. *Nature*, 512(7515):423–426, 2014.
- Darjan Salaj, Anand Subramoney, Ceca Krausnikovic, Guillaume Bellec, Robert Legenstein, and Wolfgang Maass. Spike frequency adaptation supports network computations on temporally dispersed information. *Elife*, 10:e65459, 2021.
- Andrew M Saxe, James L McClelland, and Surya Ganguli. Exact solutions to the nonlinear dynamics of learning in deep linear neural networks. *arXiv preprint arXiv:1312.6120*, 2013.
- Andrew M Saxe, Yamini Bansal, Joel Dapello, Madhu Advani, Artemy Kolchinsky, Brendan D Tracey, and David D Cox. On the information bottleneck theory of deep learning. *Journal of Statistical Mechanics: Theory and Experiment*, 2019(12):124020, 2019.
- Benjamin Scellier and Yoshua Bengio. Equilibrium propagation: Bridging the gap between energy-based models and backpropagation. *Frontiers in computational neuroscience*, 11:24, 2017.
- Louis K Scheffer, C Shan Xu, Michal Januszewski, Zhiyuan Lu, Shin-ya Takemura, Kenneth J Hayworth, Gary B Huang, Kazunori Shinomiya, Jeremy Maitlin-Shepard, Stuart Berg, et al. A connectome and analysis of the adult drosophila central brain. *Elife*, 9:e57443, 2020.
- Friedrich Schuessler, Francesca Mastrogiuseppe, Alexis Dubreuil, Srdjan Ostojic, and Omri Barak. The interplay between randomness and structure during learning in rnns. *Advances in neural information processing systems*, 33:13352–13362, 2020.
- Friedrich Schuessler, Francesca Mastrogiuseppe, Srdjan Ostojic, and Omri Barak. Aligned and oblique dynamics in recurrent neural networks. *arXiv preprint arXiv:2307.07654*, 2023.
- Mariia Seleznova and Gitta Kutyniok. Neural tangent kernel beyond the infinite-width limit: Effects of depth and initialization. In *International Conference on Machine Learning*, pp. 19522–19560. PMLR, 2022.
- Yuxiu Shao and Srdjan Ostojic. Relating local connectivity and global dynamics in recurrent excitatory-inhibitory networks. *PLOS Computational Biology*, 19(1):e1010855, 2023.
- D Simard, L Nadeau, and H Kröger. Fastest learning in small-world neural networks. *Physics Letters A*, 336(1):8–15, 2005.

- H Francis Song, Guangyu R Yang, and Xiao-Jing Wang. Training excitatory-inhibitory recurrent neural networks for cognitive tasks: a simple and flexible framework. *PLoS computational biology*, 12(2):e1004792, 2016.
- Sen Song, Per Jesper Sjöström, Markus Reigl, Sacha Nelson, and Dmitri B Chklovskii. Highly nonrandom features of synaptic connectivity in local cortical circuits. *PLoS biology*, 3(3):e68, 2005.
- Vincent Thibeault, Antoine Allard, and Patrick Desrosiers. The low-rank hypothesis of complex systems. *Nature Physics*, pp. 1–9, 2024.
- Naftali Tishby and Noga Zaslavsky. Deep learning and the information bottleneck principle. In *2015 IEEE information theory workshop (itw)*, pp. 1–5. IEEE, 2015.
- Pauli Virtanen, Ralf Gommers, Travis E Oliphant, Matt Haberland, Tyler Reddy, David Cournapeau, Evgeni Burovski, Pearu Peterson, Warren Weckesser, Jonathan Bright, et al. Scipy 1.0: fundamental algorithms for scientific computing in python. *Nature methods*, 17(3):261–272, 2020.
- Kiran Vodrahalli, Rakesh Shivanna, Maheswaran Sathiamoorthy, Sagar Jain, and Ed H Chi. Nonlinear initialization methods for low-rank neural networks. *arXiv preprint arXiv:2202.00834*, 2022.
- Johan Winnubst, Erhan Bas, Tiago A Ferreira, Zhuhao Wu, Michael N Economo, Patrick Edson, Ben J Arthur, Christopher Bruns, Konrad Rokicki, David Schauder, et al. Reconstruction of 1,000 projection neurons reveals new cell types and organization of long-range connectivity in the mouse brain. *Cell*, 179(1):268–281, 2019.
- Chloe N Winston, Dana Mastrovito, Eric Shea-Brown, and Stefan Mihalas. Heterogeneity in neuronal dynamics is learned by gradient descent for temporal processing tasks. *Neural Computation*, 35(4):555–592, 2023.
- Blake Woodworth, Suriya Gunasekar, Jason D Lee, Edward Moroshko, Pedro Savarese, Itay Golan, Daniel Soudry, and Nathan Srebro. Kernel and rich regimes in overparametrized models. In *Conference on Learning Theory*, pp. 3635–3673. PMLR, 2020.
- Lechao Xiao, Jeffrey Pennington, and Samuel Schoenholz. Disentangling trainability and generalization in deep neural networks. In *International Conference on Machine Learning*, pp. 10462–10472. PMLR, 2020.
- Marjorie Xie, Samuel Muscinelli, Kameron Decker Harris, and Ashok Litwin-Kumar. Task-dependent optimal representations for cerebellar learning. *bioRxiv*, pp. 2022–08, 2022.
- Greg Yang. Tensor programs ii: Neural tangent kernel for any architecture. *arXiv preprint arXiv:2006.14548*, 2020.
- Guangyu Robert Yang and Manuel Molano-Mazón. Towards the next generation of recurrent network models for cognitive neuroscience. *Current opinion in neurobiology*, 70:182–192, 2021.
- Guangyu Robert Yang and Xiao-Jing Wang. Artificial neural networks for neuroscientists: a primer. *Neuron*, 107(6):1048–1070, 2020.
- Anthony M Zador. A critique of pure learning and what artificial neural networks can learn from animal brains. *Nature communications*, 10(1):1–7, 2019.
- Liqiong Zhao, Bryce Beverlin, Theoden Netoff, and Duane Q Nykamp. Synchronization from second order network connectivity statistics. *Frontiers in computational neuroscience*, 5:28, 2011.

A EXTENDED DISCUSSIONS ON RELATED WORKS

Theoretical Foundations of Neural Network Regimes and Implications for Neural Circuits:

The journey of understanding deep learning systems has borne witness to unprecedented progress in the mathematical dissection of neural network functionalities (Advani et al., 2020; Jacot et al., 2018; Pezeshki et al., 2021; Baratin et al., 2021; Alemohammad et al., 2020; Yang, 2020; Agarwala et al., 2022; Atanasov et al., 2021; Azulay et al., 2021; Emami et al., 2021). These theoretical findings, until recently confined predominantly to artificial domains, have embarked upon explorations into biological neural networks, elucidating the intricate dynamics of learning and computational properties (Bordelon & Pehlevan, 2022; Liu et al., 2022a; Braun et al., 2022; Ghosh et al., 2023). Among the vanguard of these theoretical endeavors stands the dichotomy of 'rich' and 'lazy' learning regimes. Both lead to task learning, yet they carry distinct ramifications for representation and generalization (Chizat et al., 2019; Flesch et al., 2021; Geiger et al., 2020; George et al., 2022; Ghorbani et al., 2020; Woodworth et al., 2020; Paccolat et al., 2021; Nacson et al., 2022; HaoChen et al., 2021; Flesch et al., 2023). In the 'lazy' regime, which is typically associated with large initial weights, learning predominantly centers on adjusting the readout weights. This leads to minimal alterations in the network weights and representation, while capitalizing on the expansive dimensionality provided by the hidden layer's random projections (Flesch et al., 2021). In contrast, the 'rich' regime, defined by smaller initial weights, fosters the development of highly tailored hidden unit representations specifically aligned with task demands, resulting in considerable adaptations in weights and representation. It's essential to highlight that the transition and dominance between these regimes are influenced by more than just the initial weight scale. Other factors, ranging from network width to the output gain (often referred to as the α parameter), play a pivotal role (Chizat et al., 2019; Geiger et al., 2020).

A nexus between deep learning theoretical frameworks and neuroscience has unveiled applications of the rich/lazy regimes. Previous investigations characterized neural network behaviors under distinct regimes (Bordelon & Pehlevan, 2022; Schuessler et al., 2023) and discerning which mode yields solutions mimicking empirical data (Flesch et al., 2021). It is compelling to observe that the existence of multiple learning regimes isn't an isolated phenomenon in artificial systems; analogous learning patterns echo in neural circuits as well. For instance, while plasticity-driven transformations might be resource-intensive, they manifest robustly during such developmental phases, followed by minor changes afterwards (Lohmann & Kessels, 2014). Building upon these findings, our research delves deeper into the precursors of these regimes. We examine how inherent factors in the brain, especially initial weight configurations, influence the inclination towards either rich or lazy learning. This understanding is crucial for assessing the applicability of regime-specific tools in neural contexts and for shedding light on the potential benefits of having both learning regimes coexist in the brain.

Interplay of Neural Learning and structure: Understanding how the brain learns using its myriad elements is a perennial quest in neuroscience. Addressing this, certain studies have unveiled biologically plausible learning rules (Lillicrap et al., 2020; Scellier & Bengio, 2017; Diederich & Opper, 1987; Hinton, 2022; Laborieux & Zenke, 2022; Greedy et al., 2022; Sacramento et al., 2018; Payeur et al., 2021; Roelfsema & Holtmaat, 2018; Meulemans et al., 2022; Murray, 2019; Bellec et al., 2020; Liu et al., 2021; 2022b; Marschall et al., 2020), suggesting potential neural algorithms involving known neural ingredients. Concurrently, given the three primary components of a neural network's design — task, learning rule, and architecture — another avenue of research delves deep into the architectural facet, specifically focusing on how it interacts with the learning rule to enhance learning (Richards et al., 2019; Zador, 2019; Yang & Molano-Mazón, 2021). Under the structural umbrella, the neural unit's complexity and initial connectivity patterns are two crucial aspects. Complex neuron models, for instance, have shown the potential in boosting learning performance by allowing implicit forms of memory and computations at the single neuron level (Salaj et al., 2021; Winston et al., 2023). Moreover, A large body of work has investigated the effect of different random initializations on learning in deep networks (Saxe et al., 2013; Bahri et al., 2020; Glorot & Bengio, 2010; He et al., 2015; Arora et al., 2019). For instance, the variance in random initial weights can induce pronounced shifts in network behavior, ranging from the "lazy" to the "rich" regimes (Chizat et al., 2019; Flesch et al., 2021). This introduces unique inductive biases during the learning process, with distinct preferences for learning certain features (George et al., 2022). Our discourse primarily orbits around connectivity and its implications on learning dynamics in networks with simple rectified units. Our results sit within the purview of these regimes, with a widely adopted assumption of gradient

descent via backpropagation as the learning rule, while remaining open to encompassing a wider spectrum of rules in future explorations.

Neural circuit connectivity pattern and eigenspectrum: While the importance of initial weights on function and learning is clear, the impact of specific weight shapes, apart from weight scale, on rich or lazy dynamics remains less explored. The predominant focus in the literature has been on random initialization. Yet, neural circuit structures significantly diverge from this paradigm. Illustratively, one finds connectivity principles or patterns markedly different from what one observes with a mere random initialization (Pogodin et al., 2023), resulting in distinct neural dynamics; these connectivity principles or patterns include Dale’s law (Rajan & Abbott, 2006; Ipsen & Peterson, 2020; Harris et al., 2022), an over-representation of higher-order motifs (Dahmen et al., 2020) and cell-type-specific connectivity statistics (Aljadeff et al., 2015), to name a few. Given the prominence of low-rankedness observed in neural circuits (Song et al., 2005), our study centers on the influence of effective rank on the effective learning regime. As the next generation of connectivity data becomes available (Campagnola et al., 2022; MICrONS Consortium et al., 2021; Dorkenwald et al., 2022; Winnubst et al., 2019; Scheffer et al., 2020), future explorations will broaden the scope to other weight characteristics.

B PROOFS

B.1 PROOFS FOR MAIN TEXT THEOREM AND PROPOSITION

Notation Let $f(x) = W_2 W_1 x$ denote a two-layer linear network with N hidden units on d -dimensional inputs $x \in \mathbb{R}^d$, with weight matrices $W_1 \in \mathbb{R}^{N \times d}$ and $W_2 \in \mathbb{R}^{1 \times N}$. We consider m training inputs x_1, \dots, x_m and the corresponding data matrix $X = [x_1^T \dots x_m^T] \in \mathbb{R}^{d \times m}$; the output target is generated from a linear teacher network as $Y = \beta^T X$, where $\beta_i \sim \mathcal{N}(0, 1/d)$.

Since our goal is to investigate how the *shape* of the initial weights impacts network change, we will consider a fixed small (Froebenius) norm for these; i.e.,

$$\|W_1^{(0)}\|_F = \|W_2^{(0)}\|_F := \sigma \ll 1$$

We denote by s_1, \dots, s_d denote the singular values of $W_1^{(0)}$; they satisfy $\sum_{j=1}^d s_j^2 = \sigma^2$.

In what follows we focus on the whitened setting, where X has all its non zero singular values equal to 1. We also assume $m \geq d$ for simplicity (this assumption can easily be relaxed in our analysis), so that the whitened data assumption translates as $XX^T = I_d$.

Prior results Our analysis builds on prior results Atanasov et al. (2021) on the evolution of the NTK for two-layer linear networks trained by gradient flow of the mean square error. In the above setting, Atanasov et al. (2021) show that the final NTK $K^{(f)}$ (i.e. the asymptotic NTK as the number of iterations goes to infinity) is given by

$$K^{(f)} = \|\beta\| X^T (\hat{\beta} \hat{\beta}^T + I_d) X + O(\sigma^2). \quad (8)$$

where $\hat{\beta} := \beta / \|\beta\|$. We are interested in the expected kernel alignment over the tasks, in the small initialization regime:

$$\mathbb{E}_\beta [KA(K^{(f)}, K^{(0)})] := \mathbb{E}_\beta \left[\frac{\text{Tr}(K^{(f)} K^{(0)})}{\|K^{(f)}\|_F \|K^{(0)}\|_F} \right]. \quad (9)$$

Theorem 1. *In the above setting, when considering all possible initializations $W_1^{(0)}$ with small fixed norm σ , the expected kernel alignment $\mathbb{E}_\beta [KA]$ (defined in Eq. 9) is maximized with high-rank isotropic initialization, i.e with $W_1^{(0)}$ that has all its non-zero singular values equal in absolute value.*

Proof. Let us write $K^{(0)} = X^T M_0 X$ with $M_0 := W_1^{(0)T} W_1^{(0)} + \sigma^2 I_d$. Up to $O(\sigma^4)$ terms, the numerator in Eq. 9 takes the form

$$\begin{aligned} \text{Tr}(K^{(f)} K^{(0)}) &= \|\beta\| \text{Tr}(X^T (\hat{\beta} \hat{\beta}^T + I_d) X X^T M_0 X) \\ &\stackrel{(a)}{=} \|\beta\| \text{Tr}(X^T (\hat{\beta} \hat{\beta}^T + I_d) M_0 X) \\ &\stackrel{(b)}{=} \|\beta\| \text{Tr}((\hat{\beta} \hat{\beta}^T + I_d) M_0 X X^T) \\ &\stackrel{(a)}{=} \|\beta\| \text{Tr}((\hat{\beta} \hat{\beta}^T + I_d) M_0) \\ &\stackrel{(c)}{=} \|\beta\| (\hat{\beta}^T M_0 \hat{\beta} + \text{Tr} M_0) \end{aligned} \quad (10)$$

where (a) uses $XX^T = I_d$, (b) the cyclicity of the trace, and (c) the fact that $\hat{\beta}^T M_0 \hat{\beta}$ is a scalar.

As for the denominator in Eq. 9), we have,

$$\begin{aligned} \|K^{(0)}\|_F^2 &= \text{Tr}(K^{(0)} K^{(0)}) \\ &= \text{Tr}(X^T M_0 X X^T M_0 X) \\ &\stackrel{(a)}{=} \text{Tr}(M_0^2) \end{aligned} \quad (11)$$

and, up to $O(\sigma^4)$ terms,

$$\begin{aligned}
\|K^{(f)}\|_F^2 &= \text{Tr}(K^{(f)}K^{(f)}) \\
&= \|\beta\|^2 \text{Tr}(X^T(\hat{\beta}\hat{\beta}^T + I_d)XX^T(\hat{\beta}\hat{\beta}^T + I_d)) \\
&= \|\beta\|^2 \text{Tr}(X^T(\hat{\beta}\hat{\beta}^T + I_d)X) \\
&\stackrel{(a)}{=} \|\beta\|^2 \text{Tr}(\hat{\beta}\hat{\beta}^T + I_d)^2 \\
&\stackrel{(b)}{=} \|\beta\|^2(d+3)
\end{aligned} \tag{12}$$

where (a) in these two calculations uses $XX^T = I_d$ and the cyclicity of the trace; and (b) notes that the $d \times d$ matrix $\hat{\beta}\hat{\beta}^T + I_d$ has $d - 1$ eigenvalues equal to 1 and one equal to 2. Eq. 11 and 12 yield

$$\|K^{(f)}\|_F \|K^{(0)}\|_F = \|\beta\| \sqrt{(d+3) \text{Tr} M_0^2} \tag{13}$$

Putting together Eq. 10, 13, we obtain, up to additive $O(\sigma^2)$ terms,

$$\text{KA}(K^{(f)}, K^{(0)}) = \frac{\hat{\beta}^T M_0 \hat{\beta} + \text{Tr} M_0}{\sqrt{(d+3) \text{Tr} M_0^2}} \tag{14}$$

Next, averaging over the tasks requires computing the Gaussian average

$$A[M_0] := \mathbb{E}_\beta \left[\hat{\beta}^T M_0 \hat{\beta} \right] = \mathbb{E}_\beta \left[\frac{\beta^T M_0 \beta}{\|\beta\|^2} \right].$$

Lemma 1. *The map A is invariant under the action of the orthogonal group, i.e $A[UMU^T] = A[M]$ for all $M \in \mathbb{R}^{d \times d}$ and all orthogonal matrices $U \in \mathbb{R}^{d \times d}$.*

Proof. This is a consequence of the invariance of the Gaussian measure under the action of the orthogonal group. Explicitly, given an orthogonal matrix U ,

$$\begin{aligned}
A[UMU^T] &= \frac{1}{(2\pi d)^{d/2}} \int d^d \beta e^{-\|\beta\|^2/d} \left[\frac{\beta^T U M U^T \beta}{\|\beta\|^2} \right] \\
&\stackrel{\beta' := U^T \beta}{=} \frac{1}{(2\pi d)^{d/2}} \int d^d \beta' |\det U| e^{-\|U\beta'\|^2/d} \left[\frac{\beta'^T M \beta'}{\|U\beta'\|^2} \right] \\
&= \frac{1}{(2\pi d)^{d/2}} \int d^d \beta' e^{-\|\beta'\|^2/d} \left[\frac{\beta'^T M \beta'}{\|\beta'\|^2} \right] \\
&= A[M]
\end{aligned} \tag{15}$$

where the third equality follows from $|\det U| = 1$ and $\|U\beta\| = \|\beta\|$. \square

Lemma 2. *There is a constant c such that $A[M] = c \text{Tr}(M)$ for any symmetric matrix M .*

Proof. Given a symmetric matrix M , it can be diagonalized as $M = UDU^T$ where $D = \text{Diag}(\mu_1, \dots, \mu_d)$ is diagonal and U is orthogonal. By rotation invariance from Lemma 1, we have $A[M] = A[D]$, and

$$A[D] = \mathbb{E}_\beta \left[\hat{\beta}^T D \hat{\beta} \right] = \mathbb{E}_\beta \left[\sum_{j=1}^d \hat{\beta}_j^2 \mu_j \right] = \sum_{j=1}^d \mathbb{E}_\beta \left[\frac{\beta_j^2}{\|\beta\|^2} \right] \mu_j := \sum_{j=1}^d c_j \mu_j \tag{16}$$

We conclude by noting that, by invariance of the (isotropic) Gaussian measure under permutation of the vector components, the coefficients c_j are independent of j , i.e $c_j \equiv c$ for all j . In sum,

$$A[M] = A[D] = c \text{Tr} D = c \text{Tr} M. \tag{17}$$

\square

The expected kernel alignment thus takes the form,

$$\mathbb{E}_\beta[\text{KA}(K^{(f)}, K^{(0)})] = \frac{(1+c) \text{Tr } M_0}{\sqrt{(d+3) \text{Tr } M_0^2}} \quad (18)$$

up to additive $O(\sigma^2)$ terms. Finally, we note that

$$\begin{aligned} \text{Tr } M_0 &= \text{Tr}(W_1^{(0)T} W_1^{(0)} + \sigma^2 I_d) \\ &= \|W_1^{(0)}\|_F^2 + d\sigma^2 \\ &= (d+1)\sigma^2 \end{aligned} \quad (19)$$

and

$$\begin{aligned} \text{Tr } M_0^2 &= \text{Tr}(W_1^{(0)T} W_1^{(0)} + \sigma I_d)^2 \\ &= \sum_{j=1}^d (s_j^2 + \sigma^2)^2 \\ &= \sum_{j=1}^d s_j^4 + 2\sigma^2 \sum_{j=1}^d s_j^2 + d\sigma^4 \\ &= \sum_{j=1}^d s_j^4 + (d+2)\sigma^4 \end{aligned} \quad (20)$$

Substituting into Eq. 18, we have, up to additive $O(\sigma^2)$ terms,

$$\mathbb{E}_\beta[\text{KA}(K^{(f)}, K^{(0)})] = \frac{(1+c)(d+1)}{\sqrt{(d+3)(d+2 + \sum_{j=1}^d (s_j/\sigma)^4)}} \quad (21)$$

Finally, we see in Eq 21 that the maximization of $\mathbb{E}_\beta[\text{KA}]$ reduces to the following convex constrained optimization problem:

$$\min_s \sum_j s_j^4, \quad \text{subject to } \sum_j s_j^2 = \sigma^2. \quad (22)$$

The KKT solutions satisfy $s_i^2 = \sigma^2/d$ for all $j = 1 \cdots d$. This implies that the expected tangent kernel alignment is maximized when the initial weight singular values $|s_i|$ are distributed evenly across dimensions, which corresponds to a high-rank initialization. \square

Proposition 1. *Following the setup and assumptions in Theorem 1, rank-1 initialization with $W_1^{(0)} = \sigma[\hat{\beta}^T \quad \vec{0} \quad \dots \quad \vec{0}]$ leads to maximal alignment, i.e, $\text{KA}(K^{(f)}, K^{(0)}) = 1$ up to additive $O(\sigma^2)$ terms.*

Proof. We indeed have,

$$\begin{aligned} K^{(0)} &= X^T (W_1^{(0)T} W_1^{(0)} + \|W_2^{(0)}\|^2 I) X \\ &= \sigma^2 X^T (\hat{\beta} \hat{\beta}^T + I) X \end{aligned} \quad (23)$$

Thus, writing $K := X^T (\hat{\beta} \hat{\beta}^T + I) X$ and using Eq. 8, the alignment takes the form

$$\begin{aligned} \text{KA}(K^{(f)}, K^{(0)}) &:= \frac{\text{Tr}(K^{(f)} K^{(0)})}{\|K^{(f)}\|_F \|K^{(0)}\|_F} \\ &= \frac{\text{Tr}(K(K + O(\sigma^2)))}{\|K\|_F \|K + O(\sigma^2)\|_F} \\ &= \frac{\text{Tr}(K^2)}{\|K\|_F^2} + O(\sigma^2) \\ &= 1 + O(\sigma^2) \end{aligned} \quad (24)$$

\square

B.2 LEARNING REQUIREMENT BASED ON $W_h^{(0)}$ RANK

The focus of this idea is to show that no changes to hidden weights W_h is not possible (e.g. reservoir settings) for zero-error when the initial weight rank falls below a certain threshold. Freezing the hidden weights W_h would be a special case of lazy learning.

Proposition 2. Consider a linear RNN with input at time t as $X_t \in \mathbb{R}^{N \times d}$ (for $t = 1, \dots, T - 1$), target output $Y \in \mathbb{R}^{N_{out} \times d}$ only at the last step, recurrent weight matrix $W_h \in \mathbb{R}^{N \times N}$ and readout weight matrix $w \in \mathbb{R}^{N_{out} \times N}$. Here, N, N_{out}, d and T are the number of hidden units, number of classes, number of data points and number of time steps, respectively, and we assume $N, d > N_{out}$. Define initial recurrent weight $W_h^{(0)}$ and final recurrent weight $W_h^{(f)}$ that achieves zero error. Then, for arbitrary input X and target output Y , $W_h^{(f)} = W_h^{(0)}$ is not possible when $\text{rank}(W_h^{(0)}) < N_{out}$.

Proof. We have the following based on the assumption of the RNN structure, if zero-error learning is achieved:

$$Y = w^{(f)} W_h^{(f)} \left(\sum_{t=1}^{T-1} W_h^{(f)T-t-1} X_t \right). \quad (25)$$

We can prove by contradiction. Suppose $W_h^{(f)} = W_h^{(0)}$, then

$$Y = w^{(f)} W_h^{(0)} \left(\sum_{t=1}^{T-1} W_h^{(0)T-t-1} X_t \right). \quad (26)$$

Since Y is arbitrary, we can have $\text{rank}(Y) = N_{out}$ (by the assumption of $N, d > N_{out}$). Applying $\text{rank}(W_h^{(0)}) < N_{out}$ we have

$$\begin{aligned} \text{rank}(Y) &= \text{rank}\left(w^{(f)} W_h^{(0)} \left(\sum_{t=1}^{T-1} W_h^{(0)T-t-1} X_t \right)\right) \\ &\stackrel{(a)}{\leq} \min(\text{rank}(w^{(f)}), \text{rank}(W_h^{(0)}), \left(\sum_{t=1}^{T-1} W_h^{(0)T-t-1} X_t \right)) \\ &< N_{out}, \end{aligned} \quad (27)$$

where (a) is because $\text{rank}(W_h^{(0)}) < N_{out}$ so the minimum has to be less than N_{out} . This would contradict an arbitrary Y with $\text{rank}(Y) = N_{out}$. Thus, $W_h^{(f)} = W_h^{(0)}$ cannot happen and recurrent weights have to be adjusted. \square

C SETUP AND SIMULATION DETAILS

C.1 INITIAL LOW-RANK WEIGHTS CREATION

For the null case, we initialized with random Gaussian distributions where each weight element $W_{ij} \sim \mathcal{N}(0, g^2/N)$, with an initial weight variance of g . Unless otherwise mentioned, we set $g = 1.5$ and network size $N = 300$, though we also validated across other parameter choices (see Appendix D). Input and readout weights were initialized similarly as in Yang & Wang (2020) (see their *EIRNN.ipynb* notebook).

To create low-rank weights using SVD, we generated temporary weights $\hat{W}_{ij} \sim \mathcal{N}(0, g^2/N)$. Subsequently, we applied SVD to \hat{W} and retained the top components based on the desired rank. To ensure comparisons are made across constant initial weight magnitudes, the resultant weight matrix was rescaled to match the Frobenius norm of \hat{W} .

Furthermore, we present details for experimentally-driven low-rank weights. For block-specific statistics, we followed the setup in Figure S3 of Aljadeff et al. (2015), setting parameters as $\alpha = 0.02$, $\gamma = 10$, and $1 - \epsilon = 0.8$. These parameters substantially influence the weight eigenspectrum, as depicted in Figure S3 of Aljadeff et al. (2015); we selected these values specifically to emphasize the outliers and achieve a lower effective rank. These parameters represent the fraction of hyperexcitable neurons (population 1), gain of hyperexcitable connections, and the gain of remaining connections, respectively. For the creation of a chain motif, we employed the procedure described in Section S3.10 of Dahmen et al. (2020), setting $\tau_{chn} = 0.03$ (and $\tau_{chn} = -0.1$ for over-representation or under-representation of the chain motif, respectively). Here, we set $N = 100$. These parameters were chosen to provide enough distinctions from the null case, while still ensuring stability and effective task learning. The electron microscopy (EM) connectivity (of the V1 cortical column model) is obtained from Allen Institute (2023), which includes dendritic tree reconstructions and local axonal projections for hundreds of thousands of neurons, detailing their 0.5 billion synaptic connections. From this, we selected 198 cells, focusing on fully proofread neurons closest to the midpoint between layers 2/3 and 4. Connectivity strength for each neuron is determined by summing the volume of each post-synaptic density to target cells, distinguishing between excitatory and inhibitory cell types. For instance, if cell 'a' forms 10 synapses with cell 'b', the connection strength of connection[a,b] represents the combined volume of synaptic densities at cell 'b'. Inhibitory connections are assigned a sign of -1, while excitatory ones receive +1. For the Dale's law obeying initial connectivity, balanced initialization was done following the process in Yang & Wang (2020) with 80% excitatory and 20% inhibitory neurons (see the notebook *EIRNN.ipynb*).

It is crucial to highlight that, in testing our Theorem, which examines the effect of the **initial** weight rank, all low-rank modifications are not enforced during training (although the impact of enforcing these structures could be an interesting avenue for future exploration). Weights are adjusted freely based on gradient descent learning.

C.2 TASK AND TRAINING DETAILS

Our code is accessible at https://github.com/Helena-Yuhan-Liu/BioRNN_RichLazy. We used PyTorch Version 1.10.2 (Paszke et al., 2019). Simulations were executed on a computer server with x2 20-core Intel(R) Xeon(R) CPU E5-2698 v4 at 2.20GHz, with the average task training duration being around 10 minutes. Following the procedure in George et al. (2022), which delved deeply into effective laziness metrics, we employed gradient-descent learning with the SGD optimizer. Unless mentioned otherwise, the learning rate was $3e - 3$, but we validated that our findings remain consistent across various learning rates (see Appendix D). For stopping, we trained the neurogym tasks for 10000 SGD iterations, which led to comparable terminal losses and accuracies across initializations. For the sMNIST task, we concluded our training upon reaching 97% accuracy, a criterion informed by both published results and our computational resources. We also experimented with halving and doubling the training iterations and observed similar trends. All weights — input, recurrent and readout — were trained. For statistical analysis and significance tests, we used methods in the SciPy Package (Virtanen et al., 2020).

For the neuroscience tasks, we adopted the Neurogym framework (Molano-Mazon et al., 2022). Within this paper, these tasks are denoted as "2AF", "DMS", and "CXT", mirroring Neurogym

settings: $task = 'PerceptualDecisionMaking - v0'$, $task = 'DelayMatchSample - v0'$, and $task = 'ContextDecisionMaking - v0'$, respectively. To expedite simulations and facilitate numerous runs, we operated with $dt = \tau_m = 100ms$ and abbreviated task durations: for 2AF, settings were $stimulus = 700ms$ and $decision = 100ms$; for DMS, they were $sample = 100ms$, $delay = 500ms$, $test = 100ms$, and $decision = 100ms$; for CXT, they comprised $stimulus = 200ms$, $delay = 500ms$, and $decision = 100ms$. For these three tasks, we used a batch size of 32 and trained for 10000 iterations.

Regarding the sequential MNIST task LeCun (1998), we employed a row-by-row format to hasten simulations. Inputs were delivered via $N_{in} = 28$ units, each presenting a row’s grey-scaled value, culminating in 28 steps with network predictions rendered at the final step. Training hinged on the cross-entropy loss function; targets were provided throughout training for the neuroscience tasks, as per Neurogym implementation, and targets were provided at only a trial’s conclusion for the sequential MNIST task. For this task, we used a batch size of 200 and trained for 10000 iterations.

For the student-teacher two-layer linear network simulations in Figure 4A and Figure 3, we set $N = 1000$, $d = 2$ (also found similar trends for $d = 20$ and $d = 100$), $z = Fx$, all entries of w (or β when $F = I$) to 1 and entries of X are sampled from a uniform distribution over the interval $[-2, 2]$. We used standard Normal initialization for both W_1 and W_2 with $\sigma = 0.001$. For Figure 3, F is constructed from SVD, i.e. $F = USV^T$, with U and V generated from arbitrary orthogonal matrices, and S is a diagonal matrix consisting of the singular values with the top half of the singular values set to κ and bottom half set to 1, where κ is the condition number of F . For the aligned initialization, W_1 is initialized as given in Proposition 1 with $\beta = w^T F$ (w here is illustrated in Figure 3), and the F is replaced by its rank- $(d/2)$ truncation for the partially aligned initialization case. For the MNIST task shown in Figure 4B, we used a two-layer feedforward network with a ReLU activation function. The architecture consists of an input layer with 784 units corresponding to the image pixels, a hidden layer with 300 units, and a linear readout layer with 10 output units. The weights of the hidden layer and the readout layer were initialized similarly to the input and readout weights, respectively, used in the RNN settings.

D ADDITIONAL SIMULATIONS

We perform additional simulations to show the robustness of our main trends, Low-rank initial recurrent weights lead to greater changes (or effectively richer learning) in RNNs. We show the main trends observed in Figure 1 holds also for Uniform initialization (Figure 7), soft initial weight rank (Figure 8), various network sizes (Figure 9), learning rates (Figure 10), initial weight gains (Figure 11) and Dale’s Law constraint throughout training (Figure 17), finer simulation time step dt (Figure 12) and fixing the leading initial weight eigenvalue (Figure 13). The trends in Figure 2 also applies to the DMS task (Figure 5) and the CXT task (Figure 6). Also, without the low-rankness in the shuffled EM connectivity, the impact on effective laziness also goes away (Figure 19). In Figure 4 we confirm that the results, shown in Figure 1 and predicted by Theorem 1, are also observed in a two-layer linear network setup. Again, we find that in situations where initializations are random, higher rank initialization leads to greater tangent kernel alignment than lower rank cases. We have also tracked the evolution of kernel task alignment (Figure 14) and kernel effective rank over the course of training (Figure 15).

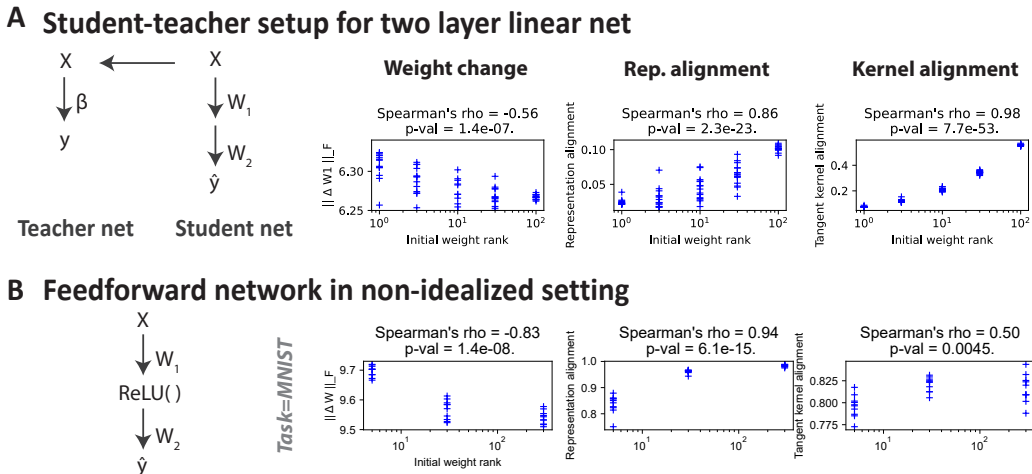


Figure 4: As predicted by the theoretical results, higher rank random initialization leads to effectively lazier learning in two-layer linear network. A) We use the student-teacher two-layer linear network setup described in Section 2.3. B) a non-idealized setting: two-layer feedforward network with ReLU activation and 300 hidden units trained on the MNIST dataset. Plotting convention follows that of Figure 1.

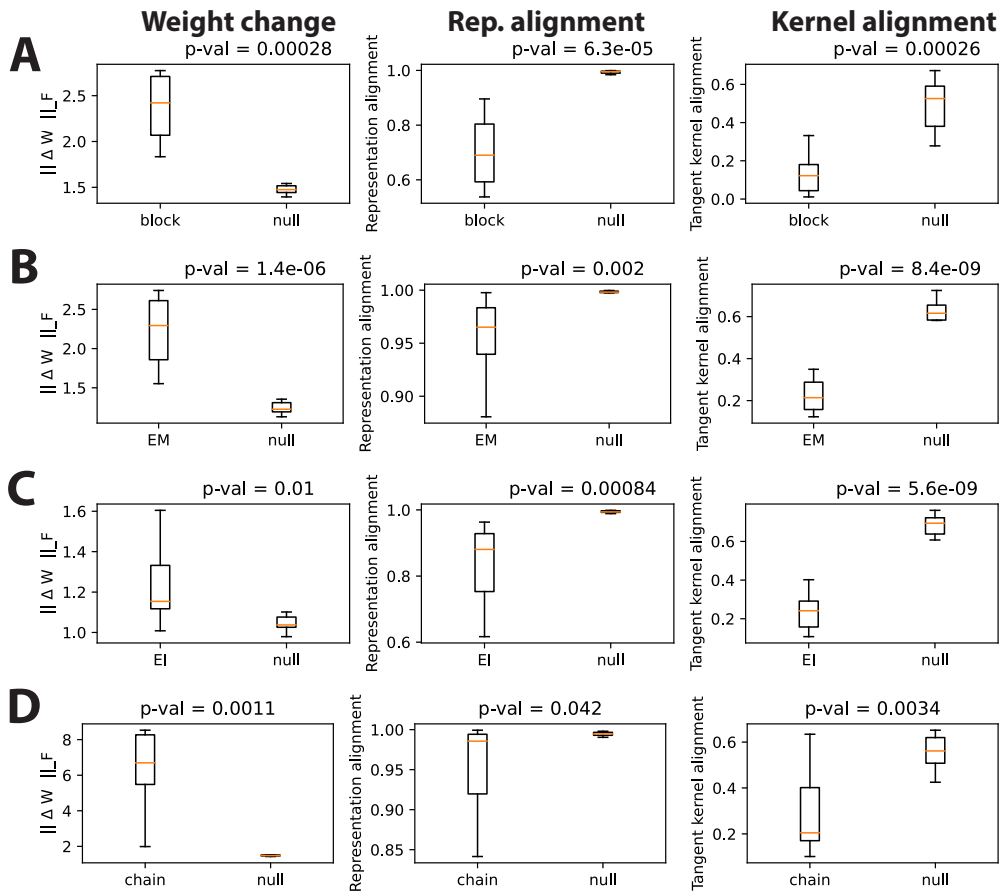


Figure 5: We repeated Figure 2 for the DMS task and observed similar trends: low-rank initialization, achieved by experimentally-driven initial connectivity in Figure 2, leads to effectively richer learning. The plotting conventions used here follow those in Figure 2, with panels A-D corresponding to the ones in that figure.

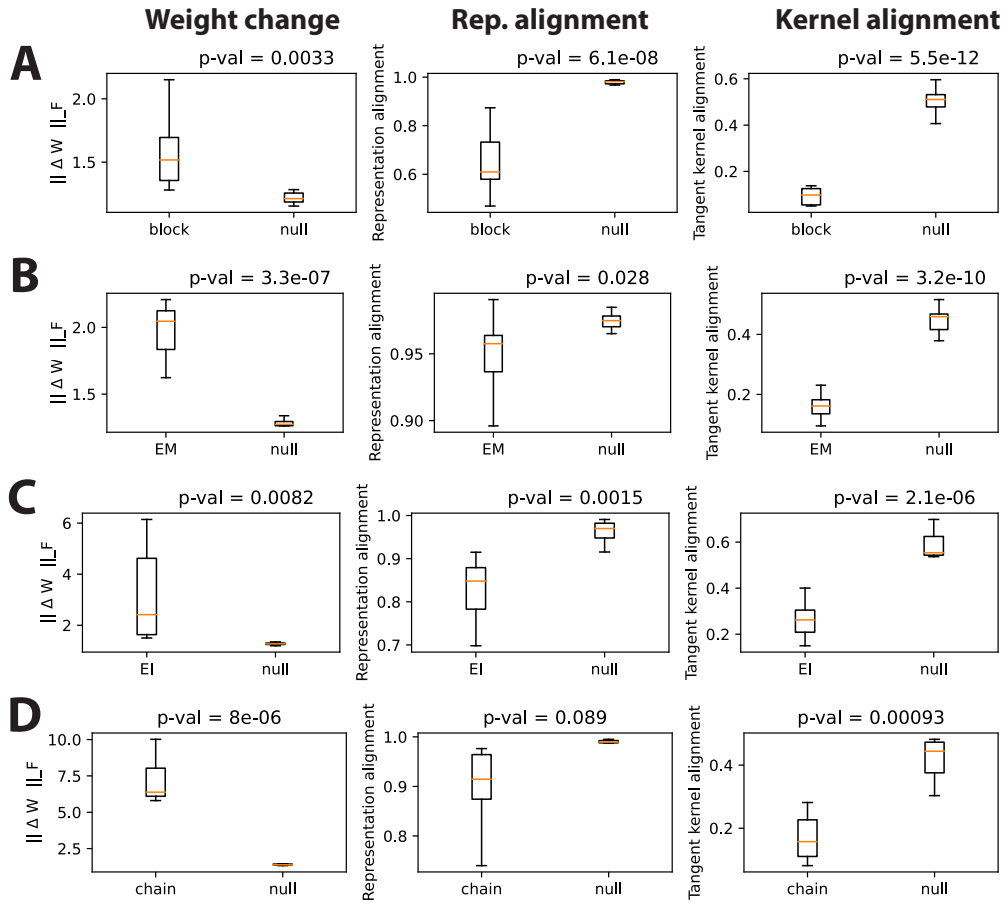


Figure 6: We repeated Figure 2 for the CXT task and observed similar trends: low-rank initialization, achieved by experimentally-driven initial connectivity in Figure 2, leads to effectively richer learning. The plotting conventions used here follow those in Figure 2, with panels A-D corresponding to the ones in that figure.

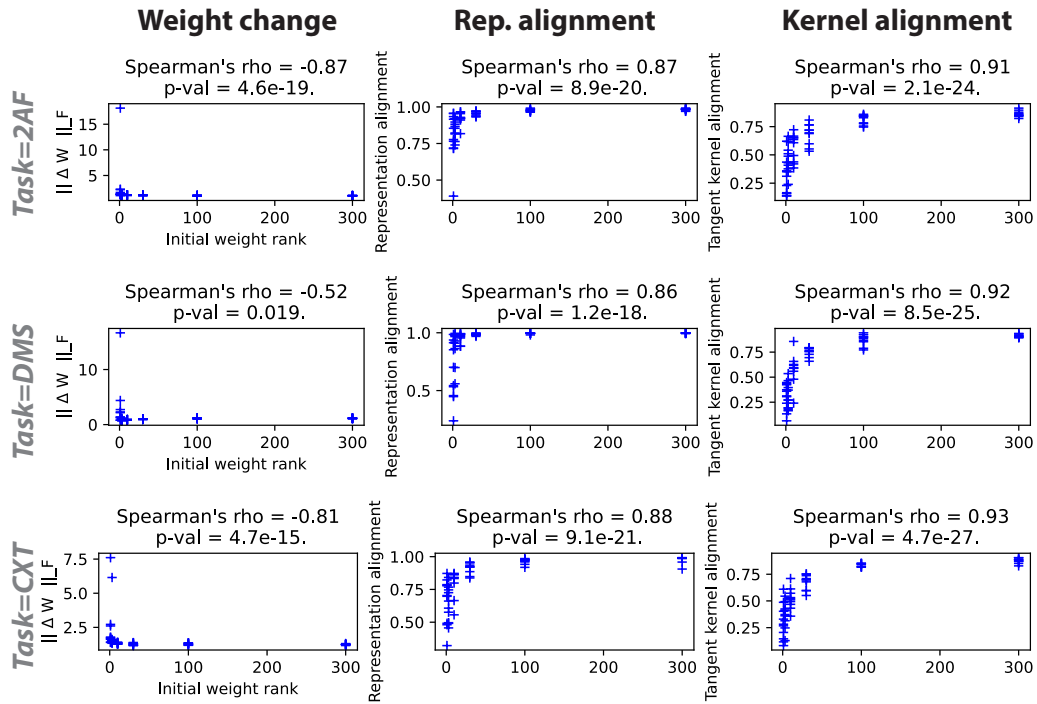


Figure 7: **Consistent trends observed in Figure 1 also for Uniform initialization.** We replicated the results of Figure 1 — where the initial weights follow a zero-mean Gaussian distribution $W_{ij} \sim \mathcal{N}(0, g^2/N)$ — but now for Uniform initialization $W_{ij} \sim \mathcal{U}\left(-\frac{g}{\sqrt{N}}, \frac{g}{\sqrt{N}}\right)$. Plotting conventions follow that of Figure 1.

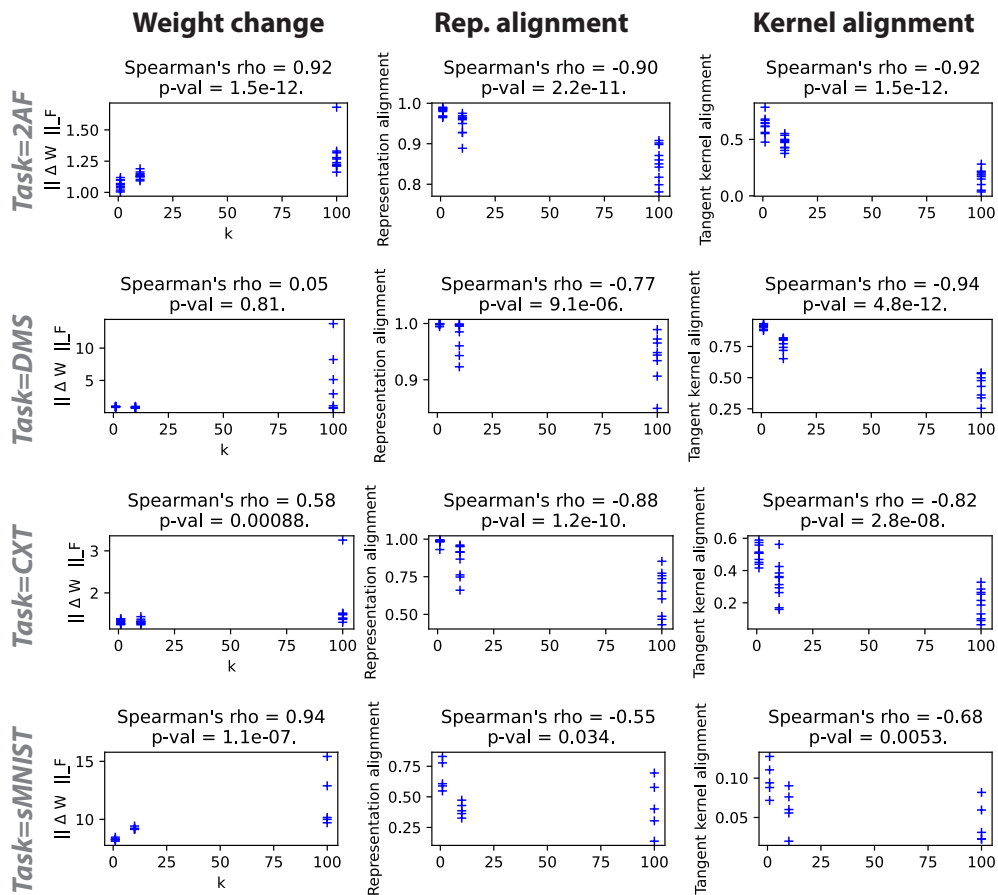


Figure 8: **Consistent trends observed in Figure 1 also for "softer" low-rank weights.** Here, instead of the "hard" low-rank weights in Figure 1 — where the i^{th} weight singular value s_i is set to 0 if $i > r$ for rank r — we introduce a smoother decay in singular value, where we replace the singular values with $s_i = s_1(1 - i/N)^k$ after performing SVD; this means that greater k leads to lower effective rank. Plotting conventions follow that of Figure 1

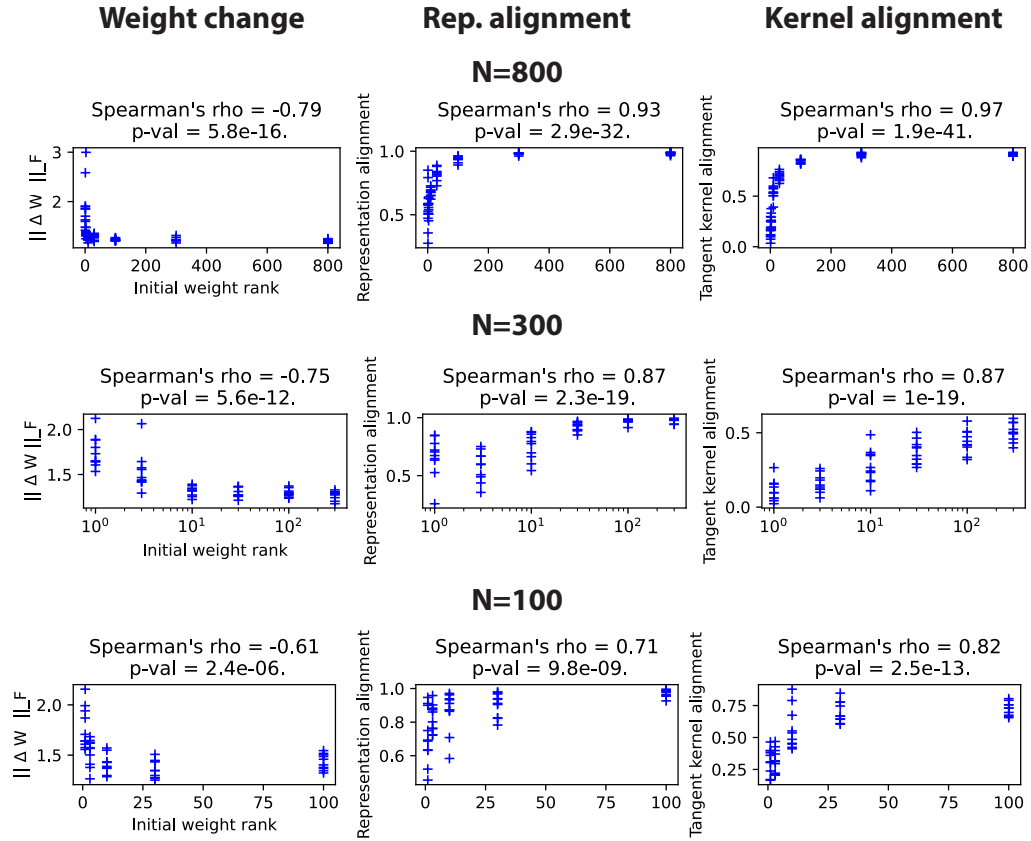


Figure 9: **Consistent trends observed in Figure 1 across various network sizes (N).** We replicated the results of Figure 1 for different values of N , using the CXT task as an illustrative example. However, the observed trend remains consistent for both the 2AF and DMS tasks. Plotting conventions follow that of Figure 1.

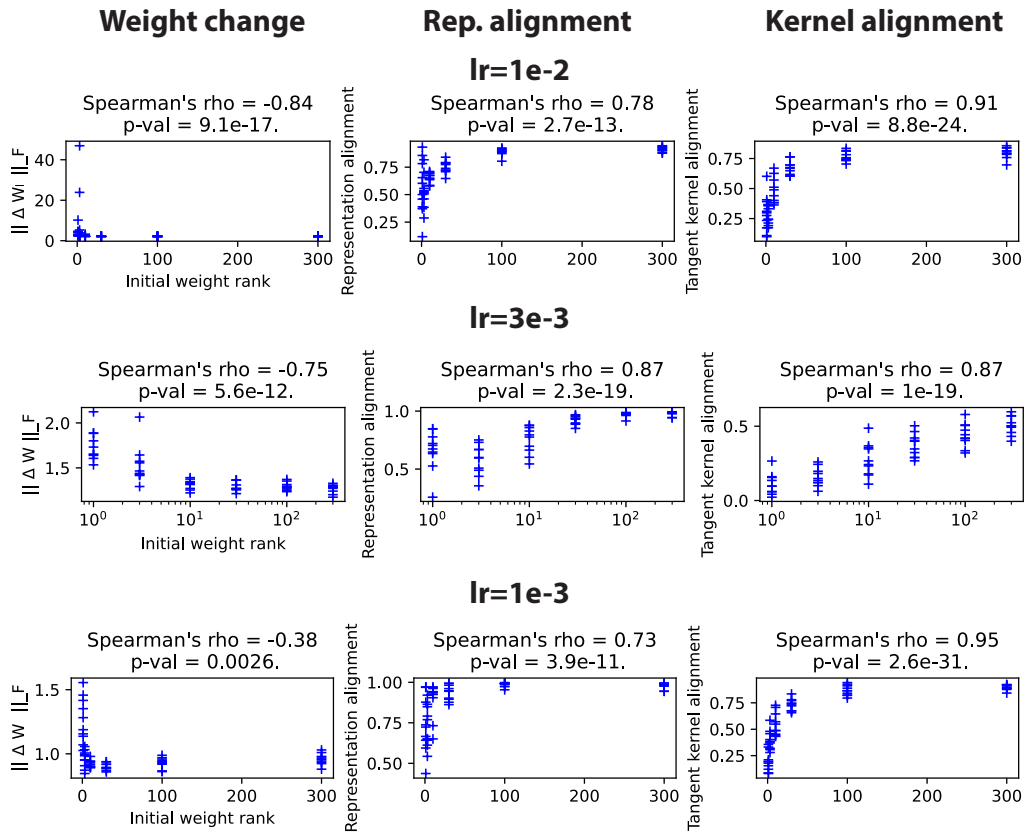


Figure 10: **Consistent trends observed in Figure 1 across various learning rates (lr).** We replicated the results of Figure 1 for different learning rates, using the CXT task as an illustrative example. However, the observed trend remains consistent for both the 2AF and DMS tasks. Plotting conventions follow that of Figure 1.

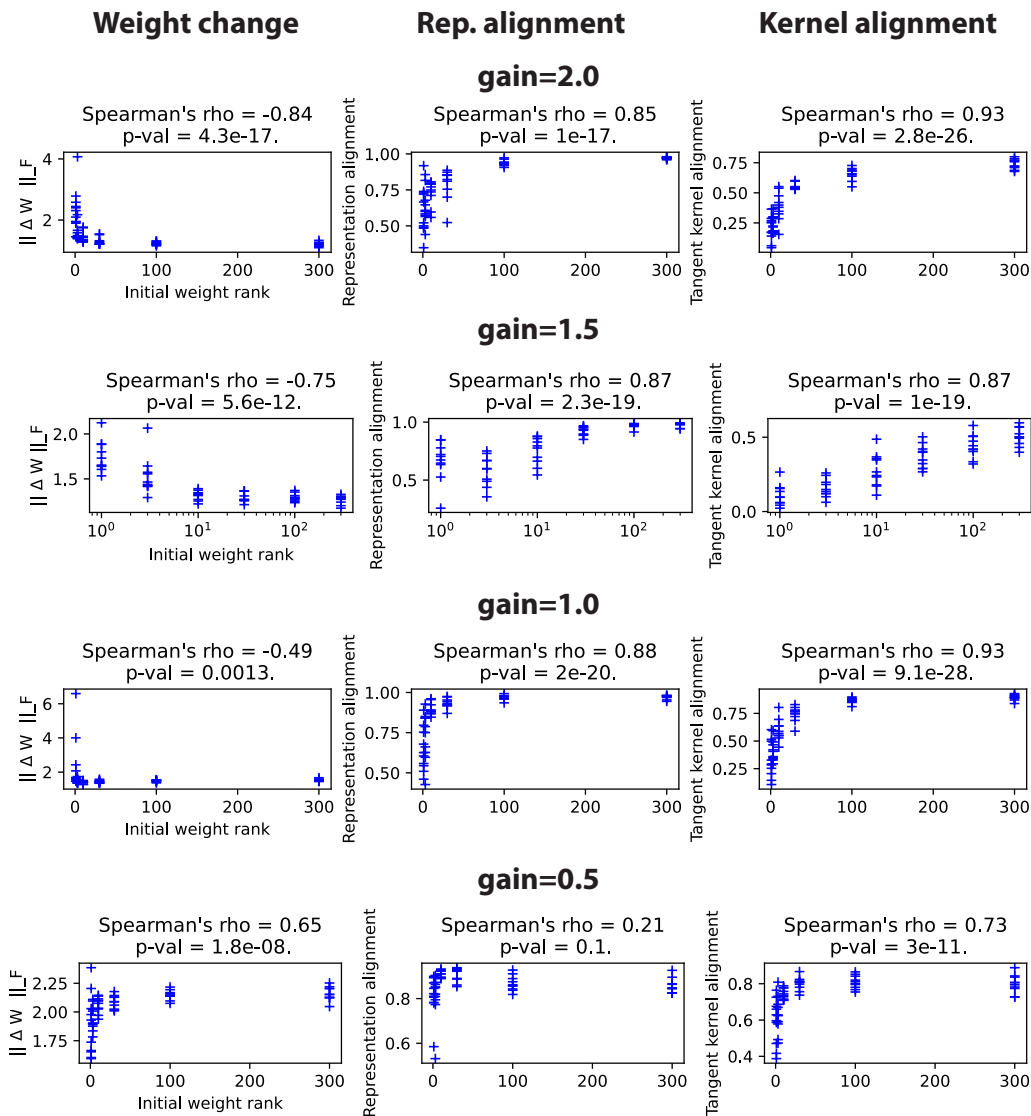


Figure 11: **Consistent trends observed in Figure 1 across various initial gain.** Here, the gain refers to g , as weights are initialized as $W_{ij} \sim \mathcal{N}(0, g^2/N)$. The trends hold for most typical range of g from 1.0 to 2.0, but gets weakened for smaller values, $g < 1.0$ (a closer examination of the regime bias in such setting in RNNs is left for future work). We replicated the results of Figure 1 for different learning rates, using the CXT task as an illustrative example. However, the observed trend remains consistent for both the 2AF and DMS tasks. Plotting conventions follow that of Figure 1.

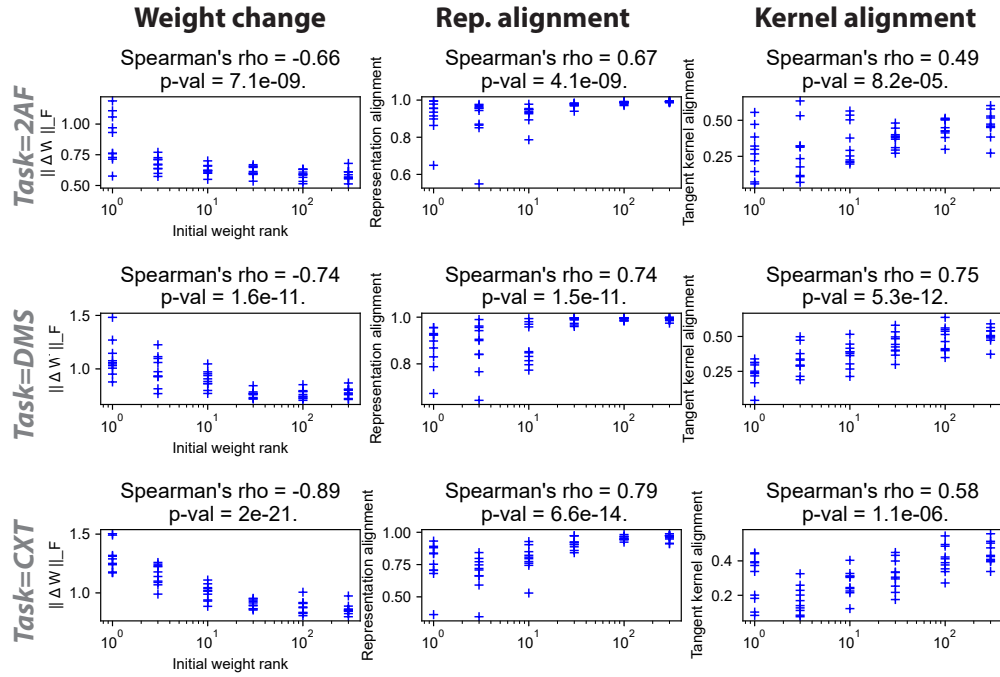


Figure 12: Trends in Figure 1 are also observed in training RNNs with a fivefold finer time step (dt) and a sequence length extended by five times. As expected, higher rank initializations led to a marked increase in effective laziness. Plotting conventions follows that of Figure 1.

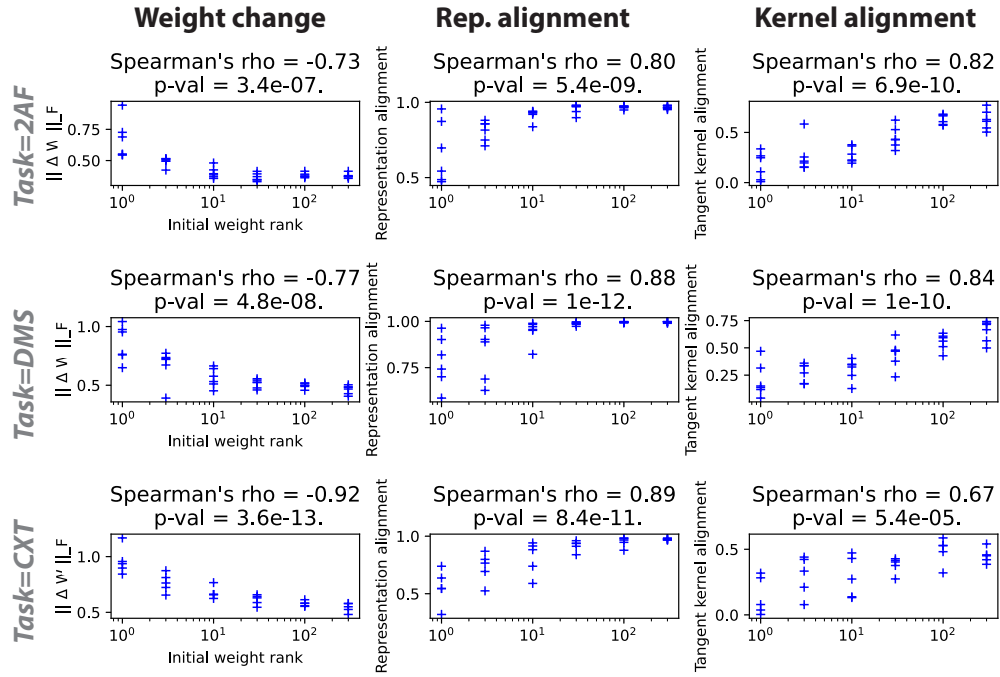


Figure 13: Trends in Figure 1 are also observed when fixing the leading weight eigenvalue instead of the Frobenius norm across comparisons. As expected, higher rank initializations lead to effectively lazier learning. Plotting conventions follows that of Figure 1.

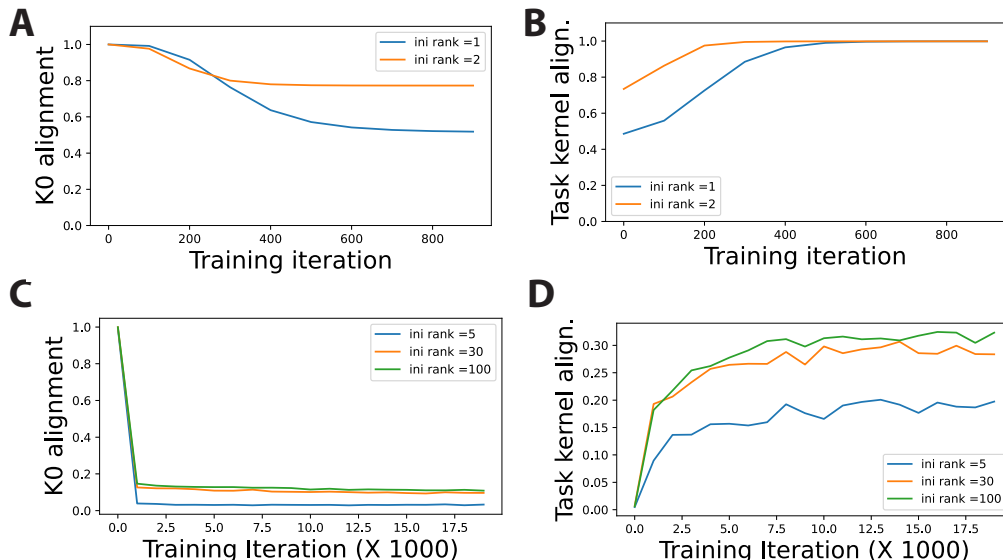


Figure 14: **[A-B] The idealized two-layer linear network setting from Fig. 2 in Atanasov et al. (2021).** A) Examining the KO alignment — the alignment between the kernel at various training iterations and the initial kernel — reveals that low-rank random initialization leads to greater changes during training; here, different curves correspond to different initial weight ranks. B) Despite these greater changes, networks with low-rank random initialization take longer to align with the task, as shown by the task kernel alignment metric $y^T K y / |y|^2 \text{Tr} K$ throughout training. We remind the reader that y corresponds to the target output and K corresponds to the NTK. **[C-D] A non-idealized setting: the sMNIST task.** C) This panel shows similar trends to A). D) Similar to B), lower-rank random initializations do not achieve as high task kernel alignment within the trained iterations. This is measured by the centered kernel alignment (CKA), which assesses the kernel’s alignment with class labels (Eq. 7 in Baratin et al. (2021)). Although higher CKA values during training could suggest enhanced feature learning (characteristic of the standard rich regime), this aligns with our findings on the effective learning regime, which focuses on changes post-training (see Introduction). Our theory in Section 2.3 suggests that lower-rank initializations require greater changes to align with the task, which would typically require more training iterations, as seen in panel B. If training is halted prematurely, perhaps due to resource constraints (as in panel D), these initializations may achieve lower final alignment within the training period. It remains unclear if extended training would lead to similar final alignment across different initializations in a wide range of scenarios. Future research should further investigate the relationship between rankedness of initializations and their impact on the converged solution’s representation, including task kernel alignment, across diverse settings.

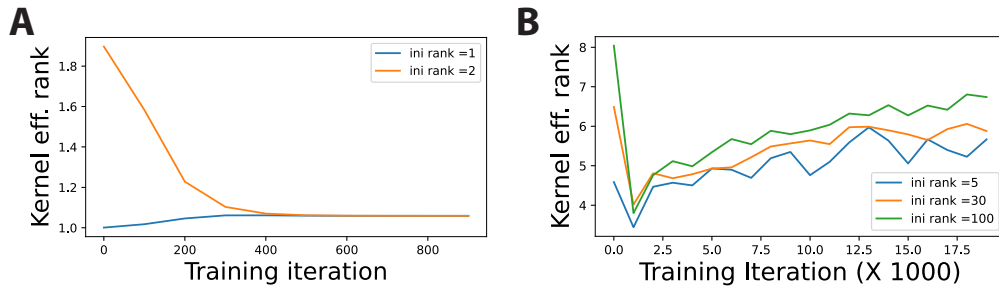


Figure 15: The evolution of the leading NTK eigenvalue relative to the rest of the eigenvalues was tracked using an effective rank measure. This measure is based on the ratio of the kernel trace to the kernel dominant eigenvalue, i.e., $\sum_i \lambda_i / \lambda_1$, which indicates the number of eigenvalues on the order of the dominant one. We apply this analysis to A) the idealized setting and B) the sMNIST task, as used in Appendix Figure 14 and Figure 1, respectively. These results suggest that the kernel effective rank approaches that of the task throughout the training process.

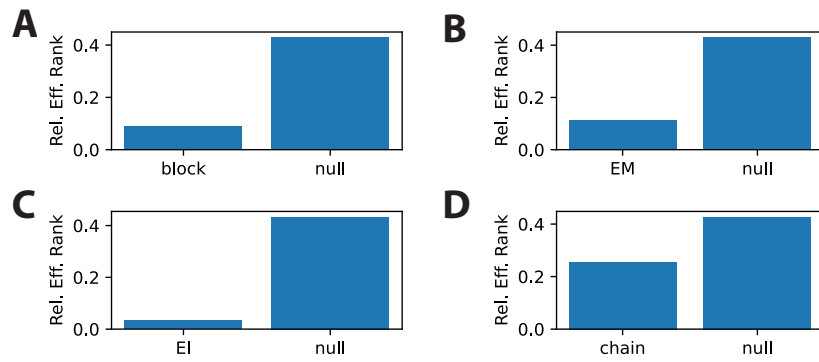


Figure 16: Measuring connectivity effective rank based on singular values instead of eigenvalues led to a similar conclusion as Figure 2: these experimentally-driven connectivity structures exhibit lower effective rank compared to random Gaussian initialization (null). The plotting conventions used here follow those in Figure 2, with panels A-D corresponding to the ones in that figure.

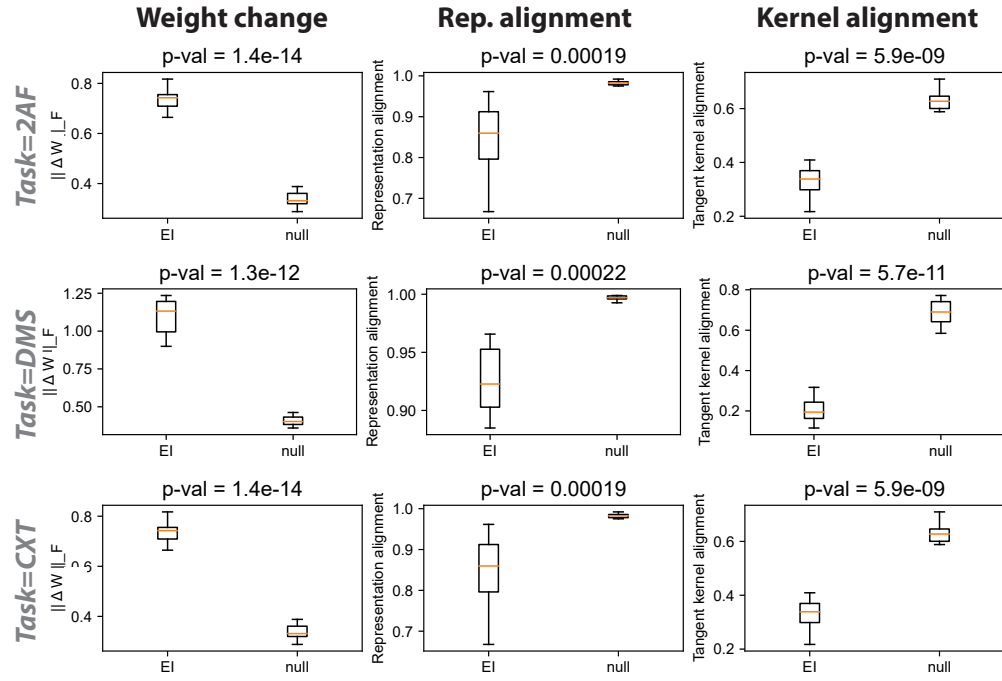


Figure 17: Maintaining the constraint of Dale’s Law during the entire training process, rather than just at initialization, produced a trend analogous to that observed in Figure 2. Plotting conventions follow that of Figure 2.

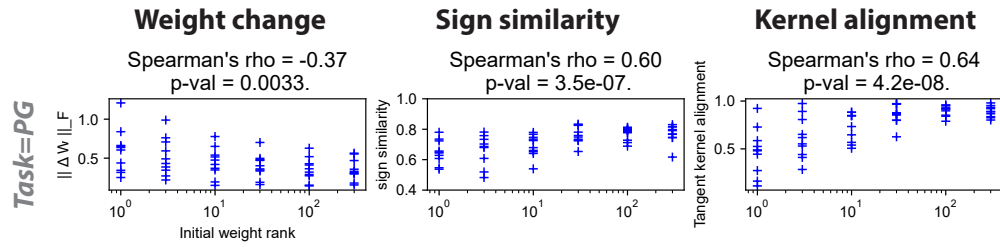


Figure 18: Training RNNs on the pattern generation task, as illustrated in Fig. S7 of Bellec et al. (2020), showed consistent trends with our conclusion: initializations with higher ranks resulted in a more pronounced tendency towards effectively lazier learning. Plotting conventions follows that of Figure 1.

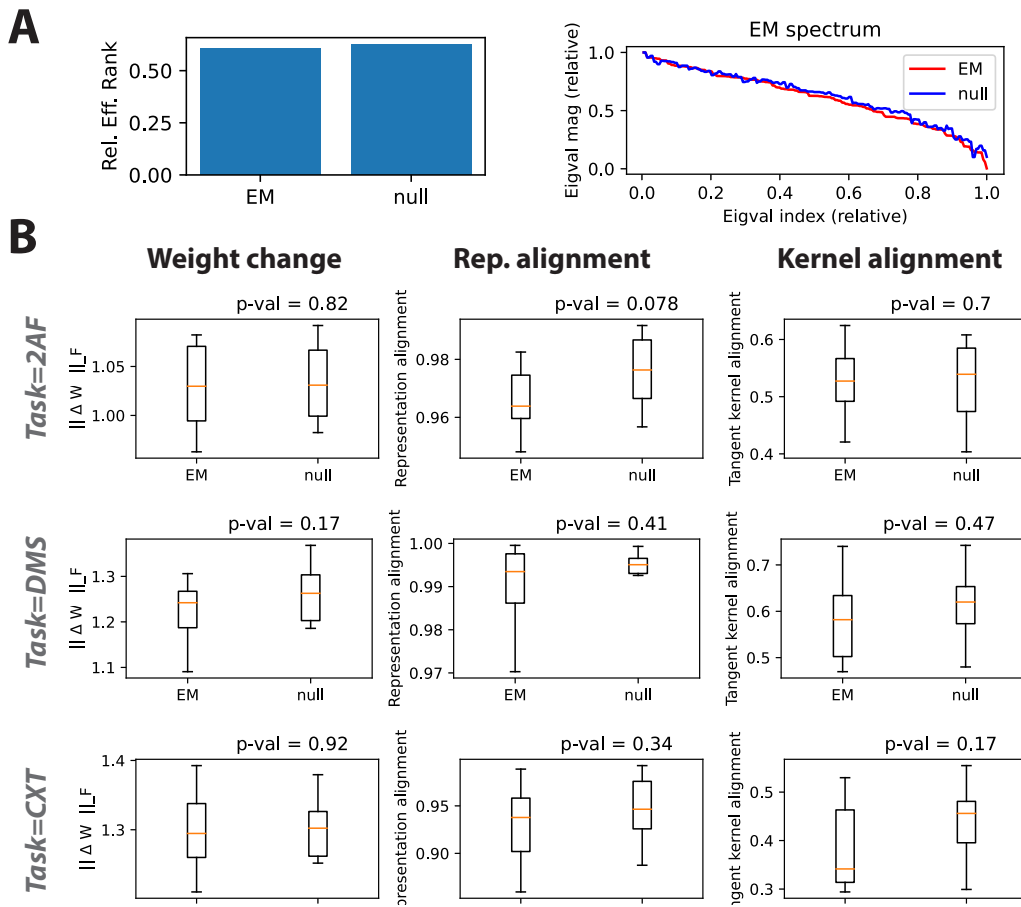


Figure 19: **Shuffling the EM connectivity, while maintaining the sparsity structure, destroys the low-rankedness and the impact on effective laziness.** We repeated the analyses with the EM initial connectivity in Figures 2 but performed random shuffling on the EM connectivity, to see if the low-rankedness and the impact on effective laziness is due to the sparsity in the dataset. Performing such shuffling destroys these trends. Plotting conventions follow that of Figure 2