
From Lazy to Rich: Exact Learning Dynamics in Deep Linear Networks

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Biological and artificial neural networks create internal representations for com-
2 plex tasks. In artificial networks, the ability to form task-specific representations
3 is shaped by datasets, architectures, initialization strategies, and optimization al-
4 gorithms. Previous studies show that different initializations lead to either a lazy
5 regime, where representations stay static, or a rich regime, where they evolve
6 dynamically. This work examines how initialization affects learning dynamics
7 in deep linear networks, deriving exact solutions for λ -balanced initializations,
8 which reflect the weight scaling across layers. These solutions explain how rep-
9 resentations and the Neural Tangent Kernel evolve from rich to lazy regimes,
10 with implications for continual, reversal, and transfer learning in neuroscience
11 and practical applications.

12 1 Introduction

13 Biological and artificial neural networks learn internal representations that enable complex tasks
14 such as categorization, reasoning, and decision-making. Both systems often develop similar repre-
15 sentations from comparable stimuli, suggesting shared information processing mechanisms Yamins
16 et al. (2014). This similarity, though not fully understood, has drawn interest from neuroscience,
17 AI, and cognitive science Haxby et al. (2001); Laakso & Cottrell (2000); Morcos et al. (2018); Ko-
18 rnblieth et al. (2019); Moschella et al. (2022). The success of neural models relies on their ability
19 to form these representations and extract relevant features from data to build internal representa-
20 tions, a complex process that in machine learning is defined by two regimes: *lazy* and *rich* Saxe
21 et al. (2014); Pennington et al. (2017); Chizat et al. (2019); Bahri et al. (2020). Despite significant
22 advances, these learning regimes and their characterization are not yet fully understood and would
23 benefit from clearer theoretical predictions, particularly regarding the influence of prior knowledge
24 (initialization) on the learning regime. We discuss related works in the appendix A.

25 **Our contributions.** (1) We derive exact solutions for the gradient flow in unequal-input-output
26 two-layer deep linear networks, under a broad range of lambda-balanced initialization conditions
27 (Section 2). (2) We model the full range of learning dynamics from *lazy* to *rich*, showing that this
28 transition is influenced by a complex interaction of architecture, *relative scale*, and *absolute scale*,
29 (Section 3). (3) We present applications relevant to both the neuroscience and machine learning
30 field, providing exact solutions for continual learning dynamics, reversal learning dynamics, and
31 transfer learning (Section 4).

32 2 Exact Learning Dynamics

33 **Preliminaries** Consider a supervised learning task where input vectors $\mathbf{x}_n \in \mathbb{R}^{N_i}$, from a set of
34 P training pairs $\{(\mathbf{x}_n, \mathbf{y}_n)\}_{n=1}^P$, need to be mapped to their corresponding target output vectors

35 $\mathbf{y}_n \in \mathbb{R}^{N_o}$. We learn this task with a two-layer linear network model that produces the output
 36 prediction $\hat{\mathbf{y}}_n = \mathbf{W}_2 \mathbf{W}_1 \mathbf{x}_n$, with weight matrices $\mathbf{W}_1 \in \mathbb{R}^{N_h \times N_i}$ and $\mathbf{W}_2 \in \mathbb{R}^{N_o \times N_h}$, where
 37 N_h is the number of hidden units. The network’s weights are optimized using full batch gradi-
 38 ent descent with learning rate η (or respectively time constant $\tau = \frac{1}{\eta}$) on the mean squared error
 39 loss $\mathcal{L}(\hat{\mathbf{y}}, \mathbf{y}) = \frac{1}{2} \langle \|\hat{\mathbf{y}} - \mathbf{y}\|^2 \rangle$, where $\langle \cdot \rangle$ denotes the average over the dataset. The dynamics are
 40 completely determined by the input covariance and input-output correlation matrices of the dataset,
 41 defined as $\tilde{\Sigma}^{xx} = \frac{1}{P} \sum_{n=1}^P \mathbf{x}_n \mathbf{x}_n^T \in \mathbb{R}^{N_i \times N_i}$ and $\tilde{\Sigma}^{yx} = \frac{1}{P} \sum_{n=1}^P \mathbf{y}_n \mathbf{x}_n^T \in \mathbb{R}^{N_o \times N_i}$, and
 42 the initialization $\mathbf{W}_2(0), \mathbf{W}_1(0)$. Our objective is to describe the entire dynamics of the network’s
 43 output and internal representations based on this initialization and the task statistics. We consider
 44 an approach first introduced in the foundational work of Fukumizu Fukumizu (1998) and extended
 45 in recent work by Braun et al. (2022), which rather than consider the dynamics of the parameters
 46 directly, we consider the dynamics of a matrix of the important statistics. In particular, defining
 47 $\mathbf{Q} = [\mathbf{W}_1 \quad \mathbf{W}_2^T]^T \in \mathbb{R}^{(N_i+N_o) \times N_h}$, we consider the $(N_i + N_o) \times (N_i + N_o)$ matrix

$$\mathbf{Q}\mathbf{Q}^T(t) = \begin{bmatrix} \mathbf{W}_1^T \mathbf{W}_1(t) & \mathbf{W}_1^T \mathbf{W}_2^T(t) \\ \mathbf{W}_2 \mathbf{W}_1(t) & \mathbf{W}_2 \mathbf{W}_2^T(t) \end{bmatrix}, \quad (1)$$

48 which is divided into four quadrants with interpretable meanings. The approach monitors sev-
 49 eral key statistics collected in the matrix. The off-diagonal blocks contain the network function
 50 $\hat{\mathbf{Y}}(t) = \mathbf{W}_2 \mathbf{W}_1(t) \mathbf{X}$, which can be used to evaluate the dynamics of the loss as shown in Fig. 1.
 51 The on-diagonal blocks capture the correlation structure of the weight matrices, allowing for the
 52 calculation of the temporal evolution of the network’s internal representations. This includes the
 53 representational similarity matrices (RSM) of the neural representations within the hidden layer, as
 54 first defined by Braun et al. (2022), $\text{RSM}_I = \mathbf{X}^T \mathbf{W}_1^T \mathbf{W}_1(t) \mathbf{X}$, $\text{RSM}_O = \mathbf{Y}^T (\mathbf{W}_2 \mathbf{W}_2^T(t))^+ \mathbf{Y}$,
 55 where $+$ denotes the pseudoinverse; and the network’s finite-width NTK Jacot et al. (2018); Lee
 56 et al. (2019); Arora et al. (2019b) $\text{NTK} = \mathbf{I}_{N_o} \otimes \mathbf{X}^T \mathbf{W}_1^T \mathbf{W}_1(t) \mathbf{X} + \mathbf{W}_2 \mathbf{W}_2^T(t) \otimes \mathbf{X}^T \mathbf{X}$, where \mathbf{I}
 57 is the identity matrix and \otimes is the Kronecker product. Hence, the dynamics of $\mathbf{Q}\mathbf{Q}^T$ describes the
 58 important aspects of network behaviour.

59 **Assumptions.** See Appendix B.2 for a further discussion of each assumptions.

- 60 • **A1 (Whitened input).** The input data is whitened, that is $\tilde{\Sigma}^{xx} = \mathbf{I}$.
- 61 • **A2 (Lambda-balanced).** The network’s weight matrices are lambda-balanced at the begin-
 62 ning of training, that is $\mathbf{W}_2(0)^T \mathbf{W}_2(0) - \mathbf{W}_1(0) \mathbf{W}_1(0)^T = \lambda \mathbf{I}$. If this condition holds at
 63 initialization, it will persist throughout training Saxe et al. (2014); Arora et al. (2018a). For
 64 completeness, we prove this in Appendix B.
- 65 • **A3 (Dimensions).** The hidden dimension of the network is defined as $N_h = \min(N_i, N_o)$,
 66 ensuring the network is neither bottlenecked ($N_h < \min(N_i, N_o)$) nor overparameterized
 67 ($N_h > \min(N_i, N_o)$).
- 68 • **A4 (Full-rank).** The input-output correlation of the task and the initial state of the network
 69 function have full rank, that is $\text{rank}(\tilde{\Sigma}^{yx}) = \text{rank}(\mathbf{W}_2(0) \mathbf{W}_1(0)) = \min(N_i, N_o)$.

70 **Lemma 2.1.** Under assumptions 1 and 2, the gradient flow dynamics of $\mathbf{Q}\mathbf{Q}^T(t)$, with initialization
 71 $\mathbf{Q}\mathbf{Q}^T(0) = \mathbf{Q}(0)\mathbf{Q}(0)^T$ can be written as a differential matrix Riccati equation

$$\tau \frac{d}{dt} (\mathbf{Q}\mathbf{Q}^T) = \mathbf{F}\mathbf{Q}\mathbf{Q}^T + \mathbf{Q}\mathbf{Q}^T \mathbf{F} - (\mathbf{Q}\mathbf{Q}^T)^2, \quad \text{where } \mathbf{F} = \begin{pmatrix} -\frac{\lambda}{2} \mathbf{I}_{N_i} & (\tilde{\Sigma}^{yx})^T \\ \tilde{\Sigma}^{yx} & \frac{\lambda}{2} \mathbf{I}_{N_o} \end{pmatrix}. \quad (2)$$

72 As derived in Fukumizu (1998) and extended in Braun et al. (2022), whenever \mathbf{F} is symmetric and
 73 diagonalizable such that $\mathbf{F} = \mathbf{P}\mathbf{\Lambda}\mathbf{P}^T$, where \mathbf{P} is an orthonormal matrix and $\mathbf{\Lambda}$ is a diagonal
 74 matrix, then the unique solution to this matrix Riccati is given by,

$$\mathbf{Q}\mathbf{Q}^T(t) = e^{\frac{\mathbf{F}t}{\tau}} \mathbf{Q}(0) \left[\mathbf{I} + \mathbf{Q}(0)^T \mathbf{P} \left(\frac{e^{2\mathbf{\Lambda} \frac{t}{\tau}} - \mathbf{I}}{2\mathbf{\Lambda}} \right) \mathbf{P}^T \mathbf{Q}(0) \right]^{-1} \mathbf{Q}(0)^T e^{\frac{\mathbf{F}t}{\tau}}. \quad (3)$$

75 In Appendix C.2 we prove that this equation is the unique solution to the initial value problem
 76 derived in Lemma 2.1 no matter the value of $\mathbf{\Lambda}$. However, as discussed in Braun et al. (2022), the
 77 solution in this form is not very useable or interpretable due to the matrix inverse mixing the blocks
 78 of $\mathbf{Q}\mathbf{Q}^T$. Additionally, we need to diagonalize \mathbf{F} . To do so we consider the compact singular value
 79 decomposition $\text{SVD}(\tilde{\Sigma}^{yx}) = \tilde{\mathbf{U}} \tilde{\mathbf{S}} \tilde{\mathbf{V}}^T$. Here, $\tilde{\mathbf{U}} \in \mathbb{R}^{N_o \times N_h}$ denote the left singular vectors, $\tilde{\mathbf{S}} \in$

80 $\mathbb{R}^{N_h \times N_h}$ the square matrix with ordered, non-zero eigenvalues on its diagonal, and $\tilde{\mathbf{V}} \in \mathbb{R}^{N_i \times N_h}$
 81 the corresponding right singular vectors. For unequal input-output dimensions ($N_i \neq N_o$), the right
 82 and left singular vectors are not square. Accordingly, for the case $N_i > N_h = N_o$, we define
 83 $\tilde{\mathbf{U}}^\perp \in \mathbb{R}^{N_o \times |N_o - N_i|}$ as a matrix containing orthogonal column vectors that complete the basis for
 84 $\tilde{\mathbf{U}}$, i.e., make $[\tilde{\mathbf{U}} \tilde{\mathbf{U}}^\perp]$ orthonormal, and $\tilde{\mathbf{V}}^\perp \in \mathbb{R}^{N_i \times |N_o - N_i|}$ as a matrix of zeros. Conversely,
 85 when $N_i = N_h < N_o$, then $\tilde{\mathbf{V}}^\perp$ is a matrix containing orthogonal column vectors that complete
 86 the basis for $\tilde{\mathbf{V}}$ and $\tilde{\mathbf{U}}^\perp$ is a matrix of zeros. Using this SVD structure we can now describe the
 87 eigendecomposition of \mathbf{F} .

88 **Lemma 2.2.** *Under assumptions 3 and 4, the eigendecomposition of $\mathbf{F} = \mathbf{P}\mathbf{\Lambda}\mathbf{P}^T$ is*

$$\mathbf{P} = \frac{1}{\sqrt{2}} \begin{pmatrix} \tilde{\mathbf{V}}(\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}}) & \tilde{\mathbf{V}}(\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}}) & \sqrt{2}\tilde{\mathbf{V}}_\perp \\ \tilde{\mathbf{U}}(\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}}) & -\tilde{\mathbf{U}}(\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}}) & \sqrt{2}\tilde{\mathbf{U}}_\perp \end{pmatrix}, \quad \mathbf{\Lambda} = \begin{pmatrix} \tilde{\mathbf{S}}_\lambda & 0 & 0 \\ 0 & -\tilde{\mathbf{S}}_\lambda & 0 \\ 0 & 0 & \lambda_\perp \end{pmatrix}, \quad (4)$$

89 where the matrices $\tilde{\mathbf{S}}_\lambda$, λ_\perp , $\tilde{\mathbf{H}}$, and $\tilde{\mathbf{G}}$ are diagonal matrices defined as:

$$\tilde{\mathbf{S}}_\lambda = \sqrt{\tilde{\mathbf{S}}^2 + \frac{\lambda^2}{4}} \mathbf{I}, \quad \lambda_\perp = \text{sgn}(N_o - N_i) \frac{\lambda}{2} \mathbf{I}_{|N_o - N_i|}, \quad \tilde{\mathbf{H}} = \text{sgn}(\lambda) \sqrt{\frac{\tilde{\mathbf{S}}_\lambda - \tilde{\mathbf{S}}}{\tilde{\mathbf{S}}_\lambda + \tilde{\mathbf{S}}}}, \quad \tilde{\mathbf{G}} = \frac{1}{\sqrt{\mathbf{I} + \tilde{\mathbf{H}}^2}}. \quad (5)$$

90 **Main theorem.** Thanks to the eigendecomposition of \mathbf{F} we can separate the solution provided in
 91 equation 3 into four quadrants. Following an approach used in Braun et al. (2022), we will find it
 92 useful to define the following variables of the initialization that will allow us to define the product
 93 $\mathbf{P}^T \mathbf{Q}(0)$ more succinctly,

$$\mathbf{B} = \mathbf{W}_2(0)^T \tilde{\mathbf{U}}(\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}}) + \mathbf{W}_1(0)^T \tilde{\mathbf{V}}(\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}}) \in \mathbb{R}^{N_h \times N_h}, \quad (6)$$

$$\mathbf{C} = \mathbf{W}_2(0)^T \tilde{\mathbf{U}}(\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}}) - \mathbf{W}_1(0)^T \tilde{\mathbf{V}}(\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}}) \in \mathbb{R}^{N_h \times N_h}, \quad (7)$$

$$\mathbf{D} = \mathbf{W}_2(0)^T \tilde{\mathbf{U}}_\perp + \mathbf{W}_1(0)^T \tilde{\mathbf{V}}_\perp \in \mathbb{R}^{N_h \times |N_o - N_i|}. \quad (8)$$

94 Using these variables of the initialization, this brings us to our main theorem:

95 **Theorem 2.3.** *Under the assumptions of whitened inputs, 1, lambda-balanced weights 2, no bottle-
 96 neck 3, and full rank 4, the temporal dynamics of $\mathbf{Q}\mathbf{Q}^T$ are*

$$\mathbf{Q}\mathbf{Q}^T(t) = \begin{pmatrix} \mathbf{Z}_1(t)\mathbf{A}^{-1}(t)\mathbf{Z}_1^T(t) & \mathbf{Z}_1(t)\mathbf{A}^{-1}(t)\mathbf{Z}_2^T(t) \\ \mathbf{Z}_2(t)\mathbf{A}^{-1}(t)\mathbf{Z}_1^T(t) & \mathbf{Z}_2(t)\mathbf{A}^{-1}(t)\mathbf{Z}_2^T(t) \end{pmatrix},$$

97 with the time-dependent variables $\mathbf{Z}_1(t) \in \mathbb{R}^{N_i \times N_h}$, $\mathbf{Z}_2(t) \in \mathbb{R}^{N_o \times N_h}$, and $\mathbf{A}(t) \in \mathbb{R}^{N_h \times N_h}$:

$$\mathbf{Z}_1(t) = \frac{1}{2} \tilde{\mathbf{V}}(\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}}) e^{\tilde{\mathbf{S}}_\lambda \frac{t}{\tau}} \mathbf{B}^T - \frac{1}{2} \tilde{\mathbf{V}}(\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}}) e^{-\tilde{\mathbf{S}}_\lambda \frac{t}{\tau}} \mathbf{C}^T + \tilde{\mathbf{V}}_\perp e^{\lambda_\perp \frac{t}{\tau}} \mathbf{D}^T, \quad (9)$$

$$\mathbf{Z}_2(t) = \frac{1}{2} \tilde{\mathbf{U}}(\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}}) e^{\tilde{\mathbf{S}}_\lambda \frac{t}{\tau}} \mathbf{B}^T + \frac{1}{2} \tilde{\mathbf{U}}(\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}}) e^{-\tilde{\mathbf{S}}_\lambda \frac{t}{\tau}} \mathbf{C}^T + \tilde{\mathbf{U}}_\perp e^{\lambda_\perp \frac{t}{\tau}} \mathbf{D}^T, \quad (10)$$

$$\mathbf{A}(t) = \mathbf{I} + \mathbf{B} \left(\frac{e^{2\tilde{\mathbf{S}}_\lambda \frac{t}{\tau}} - \mathbf{I}}{4\tilde{\mathbf{S}}_\lambda} \right) \mathbf{B}^T - \mathbf{C} \left(\frac{e^{-2\tilde{\mathbf{S}}_\lambda \frac{t}{\tau}} - \mathbf{I}}{4\tilde{\mathbf{S}}_\lambda} \right) \mathbf{C}^T + \mathbf{D} \left(\frac{e^{\lambda_\perp \frac{t}{\tau}} - \mathbf{I}}{\lambda_\perp} \right) \mathbf{D}^T. \quad (11)$$

98 The proof of Theorem 2.3 is in Appendix C. With this solution we can calculate the exact temporal
 99 dynamics of the loss, network function, RSMs and NTK (Fig. 1A, C) over a range of lambda-
 100 balanced initializations. **Implementation and simulation.** Simulation details are in Appendix F.7.

101 3 Rich and Lazy Learning

102 In this section we use these solutions to gain a deeper understanding of the transition between the
 103 *rich* and *lazy* regimes by examining the dynamics as a function of lambda – the *relative scale* – as it
 104 varies between positive and negative infinity.

105 **Dynamics of the singular values.** Here we examine a *lambda-balanced* linear network initial-
 106 ized with *task-aligned* weights. Previous research Saxe et al. (2019a) has demonstrated that initial
 107 weights that are aligned with the task remain aligned throughout training, restricting the learning
 108 dynamics to the singular values of the network.

109 **Theorem 3.1.** *Under the assumptions of Theorem 2.3 and with a task-aligned initialization, as de-
 110 fined in Saxe et al. (2013), the network function is given by the expression $\mathbf{W}_2 \mathbf{W}_1(t) = \tilde{\mathbf{U}} \mathbf{S}(t) \tilde{\mathbf{V}}^T$*

111 where $\mathbf{S}(t) \in \mathbb{R}^{N_h \times N_h}$ is a diagonal matrix of singular values with elements $s_\alpha(t)$ that evolve ac-
 112 cording to the equation, $s_\alpha(t) = s_\alpha(0) + \gamma_\alpha(t; \lambda) (\tilde{s}_\alpha - s_\alpha(0))$, where \tilde{s}_α is the α singular value
 113 of $\tilde{\mathbf{S}}$ and $\gamma_\alpha(t; \lambda)$ is a λ -dependent monotonic transition function for each singular value that in-
 114 creases from $\gamma_\alpha(0; \lambda) = 0$ to $\lim_{t \rightarrow \infty} \gamma_\alpha(t; \lambda) = 1$ defined explicitly in Appendix D.1. We find that
 115 under different limits of λ , the transition function converges pointwise to the sigmoidal ($\lambda \rightarrow 0$) and
 116 exponential ($\lambda \rightarrow \pm\infty$) transition functions,

$$\lim_{\lambda \rightarrow 0} \gamma_\alpha(t; \lambda) \rightarrow \frac{e^{2\tilde{s}_\alpha \frac{t}{\tau}} - 1}{e^{2\tilde{s}_\alpha \frac{t}{\tau}} - 1 + \frac{\tilde{s}_\alpha}{s_\alpha(0)}}, \quad \lim_{\lambda \rightarrow \pm\infty} \gamma_\alpha(t; \lambda) \rightarrow 1 - e^{-|\lambda| \frac{t}{\tau}}. \quad (12)$$

117 The proof for Theorem 3.1 can be found
 118 in Appendix D.1. As shown in Fig.4 B,
 119 as λ approaches zero, the dynamics re-
 120 semble sigmoidal learning curves that tra-
 121 verse between saddle points, characteris-
 122 tic of the *rich* regime Braun et al. (2022).
 123 In this regime the network learns the most
 124 salient features first, which can be benefi-
 125 cial for generalization Lampinen & Gan-
 126 guli (2018). Conversely, as shown in Fig.4
 127 A and C, as the magnitude of λ increases,
 128 the dynamics become exponential, charac-
 129 teristic of the *lazy* regime. In this regime,
 130 all features are treated equally and the net-
 131 work’s dynamics resemble that of a shallow
 132 network. *relative scale* λ has in shaping the
 133 learning dynamics, from sigmoidal to ex-
 134ponential, steering the network between the
 135 *rich* and *lazy* regimes.

136 **The dynamics of the representations.** We
 137 now consider how the representations of
 138 the individual parameters \mathbf{W}_1 and \mathbf{W}_2
 139 change through training. We note that un-
 140 der lambda-balanced initializations there is
 141 simple structure which persists throughout training that allows us to recover the dynamics of the
 142 parameters up to a time-dependent orthogonal transformation from the dynamics of $\mathbf{Q}\mathbf{Q}^T(t)$.

143 The effective singular values \mathcal{S}_λ of the corresponding weights are either up-weighted or down-
 144 weighted depending on the magnitude and sign of λ , splitting the representation into two parts as
 145 shown in theorem D.1. This division is reflected in the network’s internal representations. With our
 146 solution, $\mathbf{Q}\mathbf{Q}^T(t)$, which captures the temporal dynamics of the similarity between hidden layer
 147 activations, we can analyze the network’s internal representations in relation to the task. This allows
 148 us to determine whether the network adopts a *rich* or *lazy* representation, depending on the value of
 149 λ . Assuming convergence to the global minimum, which is guaranteed when the matrix \mathbf{B} is non-
 150 singular, the internal representation satisfies $\mathbf{W}_1^T \mathbf{W}_1 = \tilde{\mathbf{V}} \tilde{\mathbf{S}}_1^2 \tilde{\mathbf{V}}^T$ and $\mathbf{W}_2 \mathbf{W}_2^T = \tilde{\mathbf{U}} \tilde{\mathbf{S}}_2^2 \tilde{\mathbf{U}}^T$ with
 151 $\mathbf{W}_2 \mathbf{W}_1 = \tilde{\mathbf{U}} \tilde{\mathbf{S}} \tilde{\mathbf{V}}^T$. Theorem D.3 in the Appendix provides a detailed proof of this limiting behav-
 152 ior. To illustrate this, we consider a hierarchical semantic learning task¹, introduced in Saxe et al.
 153 (2014); Braun et al. (2022), where living organisms are organized according to their features (Fig.
 154 2A). The representational similarity of the task’s inputs ($\tilde{\mathbf{V}} \tilde{\mathbf{S}} \tilde{\mathbf{V}}^T$) reflects this hierarchical structure
 155 (Fig.2A). Similarly, the representational similarity of the task’s target values ($\tilde{\mathbf{U}} \tilde{\mathbf{S}} \tilde{\mathbf{U}}^T$) highlights
 156 the primary groupings of items. When training a two-layer network with *relative scale* λ equal to
 157 zero and task-agnostic initialization Mishkin & Matas (2015), the input and output representational
 158 similarity matrices (Fig.2 B) match the task’s structure upon convergence. As derived in Theorem
 159 D.4 the network is guaranteed to find a *rich* solution regardless of the *absolute scale*, meaning
 160 $\mathbf{W}_1^T \mathbf{W}_1 = \tilde{\mathbf{V}} \tilde{\mathbf{S}} \tilde{\mathbf{V}}^T$ and $\mathbf{W}_2 \mathbf{W}_2^T = \tilde{\mathbf{U}} \tilde{\mathbf{S}} \tilde{\mathbf{U}}^T$, as shown in Fig. 2 C. Hence the network learns task-
 161 specific representations. We also show that as λ approaches either positive or negative infinity, the

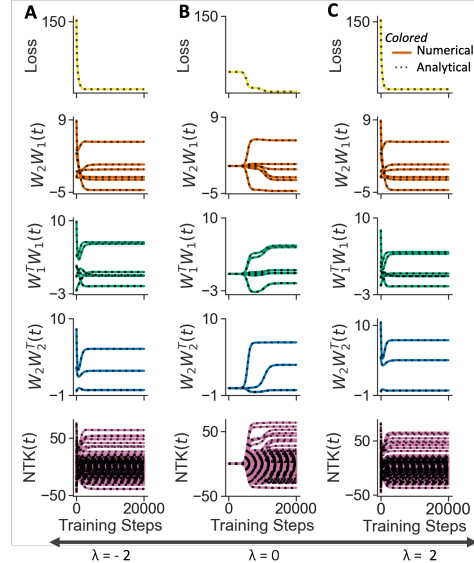


Figure 1: **A** The temporal dynamics of the numerical simulation of the loss, network dynamics function, correlation of input and output weights, and the NTK (row 1-5 respectively) are exactly matched by the analytical solution for $\lambda = -2$. **B** $\lambda = 0.001$ Large initial weight values. **C** $\lambda = 2$ initial weight values.

¹In this setting, the network has equal input and output dimensions

162 network symmetrically transitions into the *lazy* regime. As demonstrated in Theorem D.4 and illus-
 163 trated in Fig. 2, the representations converge to an identity matrix for both large positive and large
 164 negative values of λ — emerging in the output representations for large positive λ and input repre-
 165 sentations for large negative λ . This convergence indicates that the network adopts task-agnostic
 166 representations. Meanwhile, the other respective RSMs become negligible, with scales proportional
 167 to $1/\lambda$. Therefore, as shown in Theorem D.5, the NTK becomes static and equivalent to the identity
 168 matrix in the limit as λ approaches infinity. However, the downscaled representations of the net-
 169 work remain structured and task-specific. This property could be beneficial if the weights are later
 170 rescaled, such as during fine-tuning, potentially enhancing generalization and transfer learning, as
 171 we will demonstrate in Section 4. We compare this to the scenario where both weights are initial-
 172 ized with large Gaussian values, leading to *lazy* learning that maintains a fixed NTK but lacks any
 173 structural representation, as illustrated in Fig.2. Consequently, we propose a new *lazy* regime, which
 174 we refer to as the *semi-structured lazy* regime. We note that these existing regimes preserve only the
 175 input or output representation, resulting in a partial loss of structural information. All together, we
 176 find that initialization will determine which layer in the network the task specification features re-
 177 sides in: layers initialized with large values will be task-agnostic, while those initialized with small
 178 values will be task-specific.

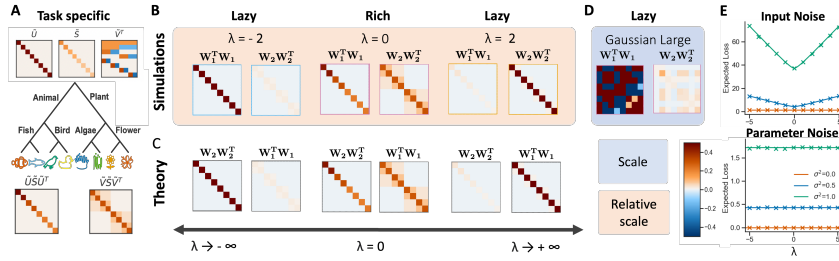


Figure 2: **A** A semantic learning task with the SVD of the input-output correlation matrix of the task. (top) U and V represent the singular vectors, and S contains the singular values. (bottom) The respective RSMs for the input and for the output task. **B** Simulation results and **C** Theoretical input and output representation matrices after training, showing convergence when initialized with varying lambda values, according to the initialization scheme described in F.7. **D** Final RSMs matrices after training converged when initialised from random large weights. **E** After convergence, the network’s sensitivity to input noise (top panel) is invariant to λ , but the sensitivity to parameter noise increases as λ becomes smaller (or larger) than zero.

179 **Representation robustness and sensitivity to noise.** Here we examine the relationship between
 180 the learning regime and the robustness of the learned representations to added noise in the inputs
 181 and parameters. The expected post-convergence loss with added noise to the inputs is determined
 182 by the norm of the network function Braun et al. (2024), which in our setting is independent of λ
 183 (Figure 2E, Appendix D.3). However, if instead noise is added to the parameters, the expected loss
 184 scales quadratically with the norm of the weight matrices Braun et al. (2024), which in our setting
 185 depend on λ . We find that under equal input-output dimensions, networks initialized with weights
 186 such that $\lambda = 0$, corresponding to the rich regime, converge to solutions that are most robust to
 187 parameter noise (Figure 2E, Appendix D.3). In practice, parameter noise could be interpreted as the
 188 noise occurring within the neurons of a biological network. Hence, a rich solution may enable a
 189 more robust representation in such systems.

190 **The impact of the architecture.** Thus far, we have found that the magnitude of the *relative scale*
 191 parameter λ determines the extent or rich and lazy learning. Here, we explore how a network’s
 192 learning regime is also shaped by the interaction of its architecture and the sign of the *relative*
 193 *scale*. We consider three types of network architectures, depicted in Fig. 3A: *funnel networks*, which
 194 narrow from input to output ($N_i > N_h = N_o$); *inverted-funnel networks*, which expand from input
 195 to output ($N_i = N_h < N_o$); and *square networks*, where input and output dimensions are equal
 196 ($N_i = N_h = N_o$). Our solution, $\mathbf{Q}\mathbf{Q}^T$, captures the dynamics of the NTK across these different
 197 network architectures. To examine the NTK’s evolution under varying λ initializations, we compute
 198 the kernel distance from initialization, as defined in Fort et al. (2020). As shown in Fig. 3B, we
 199 observe that funnel networks consistently enter the *lazy* regime as $\lambda \rightarrow \infty$, while inverted-funnel
 200 networks do so as $\lambda \rightarrow -\infty$. The NTK remains static during the initial phase, rigorously confirming
 201 the rank argument first introduced by Kunin et al. (2024) for the multi-output setting. In the opposite

202 limits of λ , these networks transition from a *lazy* regime to a *rich* regime. During this second
 203 alignment phase, the NTK matrix undergoes changes, indicating an initial *lazy* phase followed by a
 204 *delayed rich* phase. We further investigate and quantify this *delayed rich* regime, showing the NTK
 205 movement over training in Fig. 3C. This behavior is also quantified in Theorem D.6, which describes
 206 the rate of learning in this network. For square networks with equal input and output dimensions,
 207 this behavior is discussed in Section 3. Across all architectures, as $\lambda \rightarrow 0$, the networks consistently
 208 transition into the *rich* regime. Altogether, we further characterize the *delayed rich* regime in wide
 209 networks.

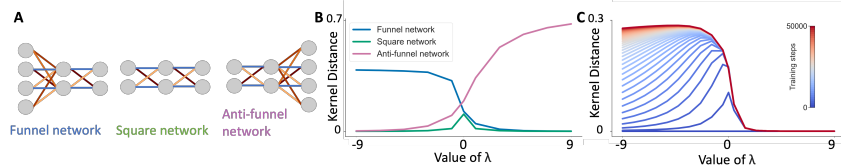


Figure 3: **A.** Schematic representations of the network architectures considered, from left to right: funnel network, square network, and inverted-funnel network. **B.** The plot shows the NTK kernel distance from initialization, as defined in Fort et al. (2020) across the three architecture depicted schematically. **C.** The NTK kernel distance away from initialization over training time.

210 4 Application

211 **Continual learning.** Similarly to the framework presented by Braun et al. (2022), our approach de-
 212 scribes the exact solutions of the networks dynamics trained across a sequence of tasks. As detailed
 213 in Appendix E.1, we demonstrate that, regardless of the chosen value of lambda, training on subse-
 214 quent tasks can result in the overwriting of previously acquired knowledge, leading to catastrophic
 215 forgetting McCloskey & Cohen (1989); Ratcliff (1990); French (1999).

216 **Reversal learning.** As demonstrated in Braun et al. (2022), reversal learning theoretically does
 217 not succeed in deep linear networks as the initialization aligns with the separatrix of a saddle point.
 218 While simulations show that the learning dynamics can escape the saddle point due to numerical
 219 imprecision, the process is catastrophically slowed in its vicinity. However, when λ is non-zero,
 220 reversal learning dynamics consistently succeed, as they avoid passing through the saddle point
 221 due to the initialization scheme. This is both theoretically proven and numerically illustrated in
 222 Appendix E.2. We also present a spectrum of reversal learning behaviors controlled by the *relative*
 223 *scale* λ , ranging from *rich* to *lazy* learning regimes. This spectrum has the potential to explain the
 224 diverse dynamics observed in animal behavior, offering insights into the learning regimes relevant
 225 to various neuroscience experiments.

226 **Transfer learning.** We consider how different λ initializations influence generalization to a new
 227 feature after being trained on an initial task. As detailed in Appendix E.3 we first train each network
 228 on the hierarchical semantic learning task described in Fig. 2. After, we add a new feature to the
 229 dataset for example ‘eats worms’ We train it specifically on the corresponding item, in this case, the
 230 goldfish, while keeping the rest of the network parameters unchanged. Afterwards, we evaluate the
 231 generalization to the other items. We observe in Appendix figure E.3 that the hierarchical structure of
 232 the data is effectively transferred to the new feature when the representation is task-specific and λ is
 233 zero. Conversely, when the output feature representation is *lazy*, meaning the hidden representation
 234 lacks adaptation, no hierarchical generalization is observed. Strikingly, when λ is positive, the
 235 hierarchical structure in the input weights remains small but structured, while the output weights
 236 exhibit a *lazy* representation and the network generalizes hierarchically. This indicates that the *lazy*
 237 regime structure can be beneficial for transfer learning.

238 5 Discussion

239 We derive exact solutions to the learning dynamics within a tractable model class: deep linear net-
 240 works. We examine the transition between the *rich* and *lazy* regimes by analyzing the dynamics
 241 as a function of λ —the *relative scale*—across its full range from positive to negative infinity. Our
 242 analysis demonstrates that the *relative scale*, λ , is pivotal in managing the transition between *rich*
 243 and *lazy* regimes.

244 References

- 245 Zeyuan Allen-Zhu, Yuanzhi Li, and Yingyu Liang. Learning and generalization in overparamete-
246 terized neural networks, going beyond two layers. *Advances in neural information processing*
247 *systems*, 32, 2019a.
- 248 Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song. A convergence theory for deep learning via over-
249 parameterization. In *International conference on machine learning*, pp. 242–252. PMLR, 2019b.
- 250 Sanjeev Arora, Nadav Cohen, Noah Golowich, and Wei Hu. A convergence analysis of gradient
251 descent for deep linear neural networks. *arXiv preprint arXiv:1810.02281*, 2018a.
- 252 Sanjeev Arora, Nadav Cohen, and Elad Hazan. On the optimization of deep networks: Implicit
253 acceleration by overparameterization. In *International Conference on Machine Learning*, pp.
254 244–253. PMLR, 2018b.
- 255 Sanjeev Arora, Nadav Cohen, Wei Hu, and Yuping Luo. Implicit regularization in deep matrix
256 factorization. *Advances in Neural Information Processing Systems*, 32, 2019a.
- 257 Sanjeev Arora, Simon S Du, Wei Hu, Zhiyuan Li, Russ R Salakhutdinov, and Ruosong Wang. On
258 exact computation with an infinitely wide neural net. *Advances in Neural Information Processing*
259 *Systems*, 32, 2019b.
- 260 Shahar Azulay, Edward Moroshko, Mor Shpigel Nacson, Blake E Woodworth, Nathan Srebro, Amir
261 Globerson, and Daniel Soudry. On the implicit bias of initialization shape: Beyond infinitesimal
262 mirror descent. In *International Conference on Machine Learning*, pp. 468–477. PMLR, 2021.
- 263 Jimmy Ba, Murat A Erdogdu, Taiji Suzuki, Zhichao Wang, Denny Wu, and Greg Yang. High-
264 dimensional asymptotics of feature learning: How one gradient step improves the representation.
265 *Advances in Neural Information Processing Systems*, 35:37932–37946, 2022.
- 266 Yasaman Bahri, Jonathan Kadmon, Jeffrey Pennington, Sam S Schoenholz, Jascha Sohl-Dickstein,
267 and Surya Ganguli. Statistical mechanics of deep learning. *Annual Review of Condensed Matter*
268 *Physics*, 11:501–528, 2020.
- 269 Pierre Baldi and Kurt Hornik. Neural networks and principal component analysis: Learning from
270 examples without local minima. *Neural networks*, 2(1):53–58, 1989.
- 271 Silvia Bernardi, Marcus K Benna, Mattia Rigotti, Jérôme Munuera, Stefano Fusi, and C Daniel
272 Salzman. The geometry of abstraction in the hippocampus and prefrontal cortex. *Cell*, 183(4):
273 954–967, 2020.
- 274 Lukas Braun, Clémentine Dominé, James Fitzgerald, and Andrew Saxe. Exact learning dynamics
275 of deep linear networks with prior knowledge. *Advances in Neural Information Processing Sys-*
276 *tems*, 35:6615–6629, 12 2022. URL [https://papers.nips.cc/paper_files/paper/2022/
277 hash/2b3bb2c95195130977a51b3bb251c40a-Abstract-Conference.html](https://papers.nips.cc/paper_files/paper/2022/hash/2b3bb2c95195130977a51b3bb251c40a-Abstract-Conference.html).
- 278 Lukas Braun, Christopher Summerfield, and Andrew Saxe. Preserving knowledge during learning
279 [unpublished manuscript]. *Department of Experimental Psychology, University of Oxford*, 2024.
- 280 Rishidev Chaudhuri, Berk Gerçek, Biraj Pandey, Adrien Peyrache, and Ila Fiete. The intrinsic
281 attractor manifold and population dynamics of a canonical cognitive circuit across waking and
282 sleep. *Nature neuroscience*, 22(9):1512–1520, 2019.
- 283 Lenaic Chizat and Francis Bach. On the global convergence of gradient descent for over-
284 parameterized models using optimal transport. *Advances in neural information processing sys-*
285 *tems*, 31, 2018.
- 286 Lenaic Chizat and Francis Bach. Implicit bias of gradient descent for wide two-layer neural networks
287 trained with the logistic loss. In *Conference on learning theory*, pp. 1305–1338. PMLR, 2020.
- 288 Lenaic Chizat, Edouard Oyallon, and Francis Bach. On lazy training in differentiable programming.
289 *Advances in neural information processing systems*, 32, 2019.

- 290 Hugo Cui, Luca Pesce, Yatin Dandi, Florent Krzakala, Yue M Lu, Lenka Zdeborová, and Bruno
291 Loureiro. Asymptotics of feature learning in two-layer networks after one gradient-step. *arXiv*
292 *preprint arXiv:2402.04980*, 2024.
- 293 Simon Du and Wei Hu. Width provably matters in optimization for deep linear neural networks. In
294 *International Conference on Machine Learning*, pp. 1655–1664. PMLR, 2019.
- 295 Simon Du, Jason Lee, Haochuan Li, Liwei Wang, and Xiyu Zhai. Gradient descent finds global
296 minima of deep neural networks. In *International conference on machine learning*, pp. 1675–
297 1685. PMLR, 2019.
- 298 Simon S Du, Xiyu Zhai, Barnabas Poczos, and Aarti Singh. Gradient descent provably optimizes
299 over-parameterized neural networks. *arXiv preprint arXiv:1810.02054*, 2018.
- 300 Matthew Farrell, Stefano Recanatesi, and Eric Shea-Brown. From lazy to rich to exclusive task rep-
301 resentations in neural networks and neural codes. *Current Opinion in Neurobiology*, 83:102780,
302 2023a. doi: 10.1016/j.conb.2023.102780.
- 303 Matthew Farrell, Stefano Recanatesi, and Eric Shea-Brown. From lazy to rich to exclusive task rep-
304 resentations in neural networks and neural codes. *Current Opinion in Neurobiology*, 83:102780–
305 102780, 12 2023b. doi: 10.1016/j.conb.2023.102780.
- 306 Timo Flesch, Keno Juechems, Tsvetomira Dumbalska, Andrew Saxe, and Christopher Summerfield.
307 Orthogonal representations for robust context-dependent task performance in brains and neural
308 networks. *Neuron*, 110:4212–4219, 12 2022. doi: 10.1016/j.neuron.2022.12.004.
- 309 Stanislav Fort, Gintare Karolina Dziugaite, Mansheej Paul, Sepideh Kharaghani, Daniel M Roy,
310 and Surya Ganguli. Deep learning versus kernel learning: an empirical study of loss landscape
311 geometry and the time evolution of the neural tangent kernel. *Advances in Neural Information*
312 *Processing Systems*, 33:5850–5861, 2020.
- 313 Robert M French. Catastrophic forgetting in connectionist networks. *Trends in cognitive sciences*,
314 3(4):128–135, 1999.
- 315 Kenji Fukumizu. Effect of batch learning in multilayer neural networks. *IEEE Transactions on*
316 *Neural Networks*, 11:17–26, 1998. doi: 10.1109/72.822506.
- 317 Mario Geiger, Stefano Spigler, Arthur Jacot, and Matthieu Wyart. Disentangling feature and lazy
318 training in deep neural networks. *Journal of Statistical Mechanics: Theory and Experiment*, 2020:
319 113301, 11 2020. doi: 10.1088/1742-5468/abc4de. URL [https://arxiv.org/pdf/1906.](https://arxiv.org/pdf/1906.08034)
320 08034.
- 321 Gauthier Gidel, Francis Bach, and Simon Lacoste-Julien. Implicit regularization of discrete gradient
322 dynamics in linear neural networks. *Advances in Neural Information Processing Systems*, 32,
323 2019.
- 324 Daniel Gissin, Shai Shalev-Shwartz, and Amit Daniely. The implicit bias of depth: How incremental
325 learning drives generalization. *arXiv preprint arXiv:1909.12051*, 2019.
- 326 James V Haxby, M Ida Gobbini, Maura L Furey, Alumit Ishai, Jennifer L Schouten, and Pietro
327 Pietrini. Distributed and overlapping representations of faces and objects in ventral temporal
328 cortex. *Science*, 293(5539):2425–2430, 2001.
- 329 Dongsung Huh. Curvature-corrected learning dynamics in deep neural networks. In *International*
330 *Conference on Machine Learning*, pp. 4552–4560. PMLR, 2020.
- 331 Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and gen-
332 eralization in neural networks. *Advances in neural information processing systems*, 31, 2018.
- 333 Arthur Jacot, François Ged, Berfin Şimşek, Clément Hongler, and Franck Gabriel. Saddle-to-saddle
334 dynamics in deep linear networks: Small initialization training, symmetry, and sparsity. *arXiv*
335 *preprint arXiv:2106.15933*, 2021.

- 336 Simon Kornblith, Mohammad Norouzi, Honglak Lee, and Geoffrey Hinton. Similarity of neural
337 network representations revisited. In *International conference on machine learning*, pp. 3519–
338 3529. PMLR, 2019.
- 339 Daniel Kunin, Atsushi Yamamura, Chao Ma, and Surya Ganguli. The asymmetric maximum margin
340 bias of quasi-homogeneous neural networks. *arXiv preprint arXiv:2210.03820*, 2022.
- 341 Daniel Kunin, Allan Raventós, Clémentine Dominé, Feng Chen, David Klindt, Andrew Saxe, and
342 Surya Ganguli. Get rich quick: exact solutions reveal how unbalanced initializations promote
343 rapid feature learning, 06 2024. URL <https://arxiv.org/abs/2406.06158>.
- 344 Aarre Laakso and Garrison Cottrell. Content and cluster analysis: assessing representational simi-
345 larity in neural systems. *Philosophical psychology*, 13(1):47–76, 2000.
- 346 Andrew K Lampinen and Surya Ganguli. An analytic theory of generalization dynamics and transfer
347 learning in deep linear networks. *arXiv preprint arXiv:1809.10374*, 2018.
- 348 Jaehoon Lee, Lechao Xiao, Samuel Schoenholz, Yasaman Bahri, Roman Novak, Jascha Sohl-
349 Dickstein, and Jeffrey Pennington. Wide neural networks of any depth evolve as linear models
350 under gradient descent. *Advances in neural information processing systems*, 32, 2019.
- 351 Sebastian Lee, Stefano Sarao Mannelli, Claudia Clopath, Sebastian Goldt, and Andrew Saxe.
352 Maslow’s hammer for catastrophic forgetting: Node re-use vs node activation. *arXiv preprint*
353 *arXiv:2205.09029*, 2022.
- 354 Aitor Lewkowycz, Yasaman Bahri, Ethan Dyer, Jascha Sohl-Dickstein, and Guy Gur-Ari. The large
355 learning rate phase of deep learning: the catapult mechanism. *arXiv preprint arXiv:2003.02218*,
356 2020.
- 357 Zhiyuan Li, Yuping Luo, and Kaifeng Lyu. Towards resolving the implicit bias of gradient descent
358 for matrix factorization: Greedy low-rank learning. *arXiv preprint arXiv:2012.09839*, 2020.
- 359 Yuhan Helena Liu, Aristide Baratin, Jonathan Cornford, Stefan Mihalas, Eric Shea-Brown, and
360 Guillaume Lajoie. How connectivity structure shapes rich and lazy learning in neural circuits.
361 *ArXiv*, 2023.
- 362 Tao Luo, Zhi-Qin John Xu, Zheng Ma, and Yaoyu Zhang. Phase diagram for two-layer relu neural
363 networks at infinite-width limit. *Journal of Machine Learning Research*, 22(71):1–47, 2021.
- 364 Sibylle Marcotte, Remi Gribonval, and Gabriel Peyré. Abide by the law
365 and follow the flow: conservation laws for gradient flows, 12 2023.
366 URL [https://proceedings.neurips.cc/paper_files/paper/2023/hash/
367 c7bee9b76be21146fd592fc2b46614d5-Abstract-Conference.html](https://proceedings.neurips.cc/paper_files/paper/2023/hash/c7bee9b76be21146fd592fc2b46614d5-Abstract-Conference.html).
- 368 Michael McCloskey and Neal J Cohen. Catastrophic interference in connectionist networks: The
369 sequential learning problem. In *Psychology of learning and motivation*, volume 24, pp. 109–165.
370 Elsevier, 1989.
- 371 Song Mei, Andrea Montanari, and Phan-Minh Nguyen. A mean field view of the landscape of two-
372 layer neural networks. *Proceedings of the National Academy of Sciences*, 115(33):E7665–E7671,
373 2018.
- 374 Dmytro Mishkin and Jiri Matas. All you need is a good init. *arXiv preprint arXiv:1511.06422*,
375 2015.
- 376 Ari Morcos, Maithra Raghu, and Samy Bengio. Insights on representational similarity in neural
377 networks with canonical correlation. *Advances in neural information processing systems*, 31,
378 2018.
- 379 Luca Moschella, Valentino Maiorca, Marco Fumero, Antonio Norelli, Francesco Locatello, and
380 Emanuele Rodolà. Relative representations enable zero-shot latent space communication. *arXiv*
381 *preprint arXiv:2209.15430*, 2022.
- 382 Srdjan Ostojic and Stefano Fusi. Computational role of structure in neural activity and connectivity.
383 *Trends in Cognitive Sciences*, 2024.

- 384 Jeffrey Pennington, Samuel Schoenholz, and Surya Ganguli. Resurrecting the sigmoid in deep learn-
385 ing through dynamical isometry: theory and practice. *Advances in neural information processing*
386 *systems*, 30, 2017.
- 387 Tomaso Poggio, Qianli Liao, Brando Miranda, Andrzej Banburski, Xavier Boix, and Jack Hidary.
388 Theory iiib: Generalization in deep networks. *arXiv preprint arXiv:1806.11379*, 2018.
- 389 David Raposo, Matthew T Kaufman, and Anne K Churchland. A category-free neural population
390 supports evolving demands during decision-making. *Nature neuroscience*, 17(12):1784–1792,
391 2014.
- 392 Roger Ratcliff. Connectionist models of recognition memory: constraints imposed by learning and
393 forgetting functions. *Psychological review*, 97(2):285, 1990.
- 394 Mattia Rigotti, Omri Barak, Melissa R Warden, Xiao-Jing Wang, Nathaniel D Daw, Earl K Miller,
395 and Stefano Fusi. The importance of mixed selectivity in complex cognitive tasks. *Nature*, 497
396 (7451):585–590, 2013.
- 397 Grant Rotskoff and Eric Vanden-Eijnden. Parameters as interacting particles: long time convergence
398 and asymptotic error scaling of neural networks. *Advances in neural information processing*
399 *systems*, 31, 2018.
- 400 Grant Rotskoff and Eric Vanden-Eijnden. Trainability and accuracy of artificial neural networks:
401 An interacting particle system approach. *Communications on Pure and Applied Mathematics*, 75
402 (9):1889–1935, 2022.
- 403 Andrew Saxe, James L. McClelland, and Surya Ganguli. Exact solutions to the nonlinear dy-
404 namics of learning in deep linear neural networks. *openreview.net*, 12 2013. URL https://openreview.net/forum?id=_wzZwKpTDF_9C.
- 405
406 Andrew M. Saxe, James L. McClelland, and Surya Ganguli. Exact solutions to the nonlinear dy-
407 namics of learning in deep linear neural networks. In *2nd International Conference on Learning*
408 *Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014, Conference Track Proceed-*
409 *ings*, 2014.
- 410 Andrew M. Saxe, James L. McClelland, and Surya Ganguli. A mathematical theory of semantic
411 development in deep neural networks. *Proceedings of the National Academy of Sciences*, 116
412 (23):11537–11546, 2019a. doi: 10.1073/pnas.1820226116. URL [https://www.pnas.org/](https://www.pnas.org/content/116/23/11537)
413 [content/116/23/11537](https://www.pnas.org/content/116/23/11537).
- 414 Andrew M Saxe, James L McClelland, and Surya Ganguli. A mathematical theory of semantic
415 development in deep neural networks. *Proceedings of the National Academy of Sciences*, 116
416 (23):11537–11546, 2019b.
- 417 Justin Sirignano and Konstantinos Spiliopoulos. Mean field analysis of neural networks: A law of
418 large numbers. *SIAM Journal on Applied Mathematics*, 80(2):725–752, 2020.
- 419 Evelyn Tang, Marcelo G Mattar, Chad Giusti, David M Lydon-Staley, Sharon L Thompson-Schill,
420 and Danielle S Bassett. Effective learning is accompanied by high-dimensional and efficient
421 representations of neural activity. *Nature neuroscience*, 22(6):1000–1009, 2019.
- 422 Salma Tarmoun, Guilherme Franca, Benjamin D Haeffele, and Rene Vidal. Understanding the
423 dynamics of gradient flow in overparameterized linear models. In *International Conference on*
424 *Machine Learning*, pp. 10153–10161. PMLR, 2021.
- 425 Zhenfeng Tu, Santiago Aranguri, and Arthur Jacot. Mixed dynamics in linear networks: Unifying
426 the lazy and active regimes. *arXiv preprint arXiv:2405.17580*, 2024.
- 427 Kay Tye, Earl Miller, Felix Taschbach, Marcus Benna, Mattia Rigotti, and Stefano Fusi. Mixed
428 selectivity: Cellular computations for complexity. *Neuron*, 112, 05 2024. doi: 10.1016/j.neuron.
429 2024.04.017.
- 430 Blake Woodworth, Suriya Gunasekar, Jason D Lee, Edward Moroshko, Pedro Savarese, Itay Golan,
431 Daniel Soudry, and Nathan Srebro. Kernel and rich regimes in overparametrized models. In
432 *Conference on Learning Theory*, pp. 3635–3673. PMLR, 2020.

- 433 Xiangxiang Xu and Lizhong Zheng. Neural feature learning in function space *. *Journal of Machine*
434 *Learning Research*, 25:1–76, 2024. URL [https://jmlr.org/papers/volume25/23-1202/](https://jmlr.org/papers/volume25/23-1202/23-1202.pdf)
435 [23-1202.pdf](https://jmlr.org/papers/volume25/23-1202/23-1202.pdf).
- 436 Yizhou Xu and Liu Ziyin. When does feature learning happen? perspective from an analytically
437 solvable model. *arXiv preprint arXiv:2401.07085*, 2024.
- 438 Daniel LK Yamins, Ha Hong, Charles F Cadieu, Ethan A Solomon, Darren Seibert, and James J
439 DiCarlo. Performance-optimized hierarchical models predict neural responses in higher visual
440 cortex. *Proceedings of the national academy of sciences*, 111(23):8619–8624, 2014.
- 441 Greg Yang and Edward J Hu. Feature learning in infinite-width neural networks. *arXiv preprint*
442 *arXiv:2011.14522*, 2020.
- 443 Greg Yang, Edward J Hu, Igor Babuschkin, Szymon Sidor, Xiaodong Liu, David Farhi, Nick Ry-
444 der, Jakub Pachocki, Weizhu Chen, and Jianfeng Gao. Tensor programs v: Tuning large neural
445 networks via zero-shot hyperparameter transfer. *arXiv preprint arXiv:2203.03466*, 2022.
- 446 Libin Zhu, Chaoyue Liu, Adityanarayanan Radhakrishnan, and Mikhail Belkin. Catapults in sgd:
447 spikes in the training loss and their impact on generalization through feature learning. *arXiv*
448 *preprint arXiv:2306.04815*, 2023.
- 449 Liu Ziyin, Botao Li, and Xiangming Meng. Exact solutions of a deep linear network. *Advances in*
450 *Neural Information Processing Systems*, 35:24446–24458, 2022.
- 451 Difan Zou, Yuan Cao, Dongruo Zhou, and Quanquan Gu. Gradient descent optimizes over-
452 parameterized deep relu networks. *Machine learning*, 109:467–492, 2020.

453 A Related Work

454 **Lazy regime.** Extensive research has identified a fundamental phenomenon in overparameterized
455 neural networks: during training, these networks frequently remain near their linearized form, un-
456 dergoing minimal changes in the parameter space Chizat et al. (2019). Consequently, they adopt
457 learning dynamics akin to kernel regression, characterized by the Neural Tangent Kernel (NTK)
458 matrix and exhibiting exponential learning behavior Du et al. (2018); Jacot et al. (2018); Du et al.
459 (2019); Allen-Zhu et al. (2019a,b); Zou et al. (2020). This behavior, known as the *lazy* or kernel
460 regime, typically occurs in infinitely wide architectures and can be triggered by large variance initial-
461 ization at the start of training Jacot et al. (2018); Chizat et al. (2019). While the *lazy* regime offers
462 valuable insights into how networks converge to a global minimum, it does not fully account for the
463 generalization capabilities of neural networks trained with standard initializations. It is, therefore,
464 widely believed that another regime, driven by small or vanishing initializations, underpins some of
465 the successes of neural networks.

466 **Rich regime.** In contrast, the *rich* feature-learning regime is characterized by a NTK that evolves
467 throughout training, accompanied by non-convex dynamics that navigate saddle points Baldi &
468 Hornik (1989); Saxe et al. (2014, 2019b); Jacot et al. (2021). This regime features sigmoidal learn-
469 ing curves and simplicity biases, such as low-rankness Li et al. (2020) or sparsity Woodworth et al.
470 (2020). Numerous studies have shown that the *absolute scale* of initialization drives the *rich* regime,
471 which typically emerges at small initialization scales Chizat et al. (2019); Geiger et al. (2020). How-
472 ever, it’s also been shown that even at small initialization scales, differences in weight magnitudes
473 between layers can induce the *lazy* learning regime Azulay et al. (2021); Kunin et al. (2024). This
474 highlights the significance of both *absolute scale* (initialization variance) and *relative scale* (differ-
475 ence in weight magnitude between layers) in generating diverse learning dynamics. Beyond *absolute*
476 *scale* and *relative scale*, additional aspects of initialization can profoundly affect feature learning,
477 including the effective rank of the weight matrices Liu et al. (2023), layer-specific initialization
478 variances Yang & Hu (2020); Luo et al. (2021); Yang et al. (2022), and the use of large learning
479 rates Lewkowycz et al. (2020); Ba et al. (2022); Zhu et al. (2023); Cui et al. (2024). These findings
480 illustrate the effect of initialization on inducing complex learning behavior through the resulting
481 dynamics. Here we develop a solvable model which captures these diverse phenomena.

482 **Rich and lazy regimes in the brain.** The distinction between *rich* and *lazy* learning may also hold
483 implications for neuroscience, where neural representations have been argued to have task-specific
484 or task-agnostic characteristics in different settings Farrell et al. (2023a); Ostojic & Fusi (2024);
485 Tye et al. (2024). The *lazy* regime can be linked to the non-linear mixed selectivity of neurons,
486 where task variables are represented in a high-dimensional space which mixes various potentially
487 relevant variables Raposo et al. (2014); Tang et al. (2019); Rigotti et al. (2013); Bernardi et al.
488 (2020). Conversely, the *rich* regime aligns with linear mixed selectivity Tye et al. (2024) and the
489 manifold learning regime, where the brain encodes tasks on a structured, low-dimensional, task-
490 specific manifold, as observed in grid cells within the entorhinal cortex Chaudhuri et al. (2019);
491 Bernardi et al. (2020); Flesch et al. (2022).

492 **Linear networks.** Our work builds upon a rich body of research on deep linear networks, which,
493 despite their simplicity, have proven to be valuable models for understanding more complex neu-
494 ral networks Baldi & Hornik (1989); Fukumizu (1998); Saxe et al. (2014). Previous research has
495 extensively analyzed convergence Arora et al. (2018a); Du & Hu (2019), generalization properties
496 Lampinen & Ganguli (2018); Poggio et al. (2018); Huh (2020), and the implicit bias of gradient
497 descent Arora et al. (2019a); Woodworth et al. (2020); Chizat & Bach (2020); Kunin et al. (2022)
498 in linear networks. These studies have also revealed that deep linear networks have intricate fixed
499 point structures and nonlinear learning dynamics in parameter and function space, reminiscent of
500 phenomena observed in nonlinear networks Arora et al. (2018b); Lampinen & Ganguli (2018).
501 Seminal work by Saxe et al. (2014) laid the groundwork by providing exact solutions to gradient
502 flow dynamics under task-aligned initializations, demonstrating that the largest singular values are
503 learned first during training. This analysis has been extended to deep linear networks Arora et al.
504 (2018b, 2019a); Ziyin et al. (2022) with more flexible initialization schemes Gidel et al. (2019);
505 Tarmoun et al. (2021); Gissin et al. (2019). This work directly builds on the matrix Riccati for-
506 mulation proposed by Fukumizu (1998); Braun et al. (2022) which extends these solutions to wide
507 networks. We extend and refine these results to obtain the dynamics for lambda-balanced initializa-
508 tion dynamics of networks to more clearly demonstrate the impact of initialization on *rich* and *lazy*
509 learning regimes also developed in Tu et al. (2024) for a set of orthogonal initializations. Our work

510 extends previous analysis Xu & Ziyin (2024); Kunin et al. (2024) of these regime to wide networks.
 511 Previous studies leveraged these solutions primarily to characterize convergence rates; however, our
 512 work goes beyond this by providing a comprehensive characterization of the complete dynamics of
 513 the system Tarmoun et al. (2021).

514 **Infinite-width networks.** Recent advances in understanding the *rich* regime have largely stemmed
 515 from examining how the initialization variance and layer-wise learning rates must scale in the
 516 infinite-width limit to maintain consistent behavior in activations, gradients, and outputs. Several
 517 studies have employed statistical mechanics tools to derive analytical solutions for the *rich* popu-
 518 lation dynamics of two-layer nonlinear neural networks initialized using the *mean field* parameter-
 519 ization Mei et al. (2018); Rotskoff & Vanden-Eijnden (2018); Chizat & Bach (2018); Sirignano &
 520 Spiliopoulos (2020); Rotskoff & Vanden-Eijnden (2022); Sirignano & Spiliopoulos (2020). Other
 521 methods for analyzing deep network dynamics include the NTK limit, where the network effectively
 522 performs kernel regression without feature learning Jacot et al. (2018); Lee et al. (2019); Arora et al.
 523 (2019b). Our solution allows us to study the evolution of the NTK and the influence of *absolute*
 524 *scale* and *relative scale* on the transition between *lazy* and *rich* learning in finite width networks
 525 Jacot et al. (2021); Xu & Ziyin (2024); Kunin et al. (2024); Chizat et al. (2019). Furthermore, these
 526 approaches typically require numerical integration or operate within a limited learning regime, and
 527 are unable to describe the learning dynamics of hidden representations. Instead, our work focuses
 528 on the impact of initialization on representation learning dynamics and derives explicit analytical
 529 solutions within tractable models.

530 B Preliminaries

531 B.1 Appendix: Balanced Condition

532 **Definition B.1** (Definition of λ -balanced property (Saxe et al. (2013), Marcotte et al. (2023))). The
 533 weights $\mathbf{W}_1, \mathbf{W}_2$ are λ -balanced if and only if there exists a **Balanced Coefficient** $\lambda \in \mathbb{R}$ such that:

$$B(\mathbf{W}_1, \mathbf{W}_2) = \mathbf{W}_2^T \mathbf{W}_2 - \mathbf{W}_1 \mathbf{W}_1^T = \lambda \mathbf{I} \quad (13)$$

534 where B is called the **Balanced Computation**.

535 For $\lambda = 0$ we have **Zero-Balanced** given as **A5** (). $\mathbf{W}_1(0) \mathbf{W}_1(0)^T = \mathbf{W}_2(0)^T \mathbf{W}_2(0)$.

536 **Theorem B.2. Balanced Condition Persists Through Training**

537 *Suppose at initialization*

$$\mathbf{W}_2(0)^T \mathbf{W}_2(0) - \mathbf{W}_1(0) \mathbf{W}_1(0)^T = \lambda \mathbf{I} \quad (14)$$

538 *Then for all $t \geq 0$*

$$\mathbf{W}_2(t)^T \mathbf{W}_2(t) - \mathbf{W}_1(t) \mathbf{W}_1(t)^T = \lambda \mathbf{I} \quad (15)$$

539 *Proof.* Consider:

$$\begin{aligned} \tau \frac{d}{dt} [\mathbf{W}_2(t) \mathbf{W}_2(t)^T - \mathbf{W}_1(t) \mathbf{W}_1(t)^T] &= \left(\tau \frac{d}{dt} \mathbf{W}_2(t) \right)^T \mathbf{W}_2(t) + \mathbf{W}_2(t)^T \left(\tau \frac{d}{dt} \mathbf{W}_2(t) \right) \\ &\quad - \left(\tau \frac{d}{dt} \mathbf{W}_1(t) \right) \mathbf{W}_1(t)^T - \mathbf{W}_1(t) \left(\tau \frac{d}{dt} \mathbf{W}_1(t) \right)^T \\ &= \mathbf{W}_1(t) \left(\tilde{\Sigma}^{yx} - \mathbf{W}_2(t) \mathbf{W}_1(t) \tilde{\Sigma}^{xx} \right)^T \mathbf{W}_2(t) \\ &\quad + \mathbf{W}_2(t)^T \left(\tilde{\Sigma}^{yx} - \mathbf{W}_2(t) \mathbf{W}_1(t) \tilde{\Sigma}^{xx} \right) \mathbf{W}_1(t) \\ &\quad - \mathbf{W}_2(t)^T \left(\tilde{\Sigma}^{yx} - \mathbf{W}_2(t) \mathbf{W}_1(t) \tilde{\Sigma}^{xx} \right) \mathbf{W}_1(t) \\ &\quad - \mathbf{W}_1(t) \left(\tilde{\Sigma}^{yx} - \mathbf{W}_2(t) \mathbf{W}_1(t) \tilde{\Sigma}^{xx} \right) \mathbf{W}_2(t) \\ &= \mathbf{0} \end{aligned}$$

540 Note that $\mathbf{W}_2(t)^T \mathbf{W}_2(t) - \mathbf{W}_1(t) \mathbf{W}_1(t)^T$ is conserved for any initial value λ . \square

541 B.2 Discussion Assumptions

542 **Whittened Inputs.** Although the whittened input assumption is quite strong, it is commonly used
 543 in analytical work to obtain exact solutions, and much of the existing literature relies on these solu-
 544 tions Fukumizu (1998); Braun et al. (2022); Kunin et al. (2024). Kunin et al. (2024) goes further by
 545 exploring the implicit bias of the trajectory without relying on exact solutions. When $X^T X$ is low-
 546 rank, they can only predict the trajectories in the limit as $\lambda \rightarrow \pm\infty$. If the interpolating manifold is
 547 one-dimensional, the solution can be solved exactly in terms of λ (black dots).

548 **Dimension.** Fukumizu assumed equal input and output dimensions $N_i = N_o$, but allowed for a
 549 bottleneck in the hidden dimension of the network $N_h \leq N_i = N_o$. The work by Braun et al. (2022)
 550 extended Fukumizu (1998) solutions to cases with unequal input and output dimensions $N_i \neq N_o$,
 551 but to so did not allow a bottleneck $N_h = \min\{N_i, N_o\}$ and added an assumption on the invertibility
 552 of a statistic of the singular vector overlap between the model and the input-output statistics. In our
 553 work we allow for unequal input and output $N_i \neq N_o$ and do not introduce an additional invertibility
 554 assumption.

555 **Balancedness.** The main distinction between our work and prior works is that both Fukumizu
 556 (1998) and Braun et al. (2022) assumed zero-balanced $\mathbf{W}_1(0) \mathbf{W}_1(0)^T = \mathbf{W}_2(0)^T \mathbf{W}_2(0)$, while
 557 we relax this assumption to λ -balanced. The zero-balanced condition restricts the networks to a
 558 *rich* setting. We develop solutions to explore the continuum between the *rich* and the *lazy* regime.
 559 While some works, such as Tarmoun et al. (2021), have considered removing this constraint, their
 560 solutions remain in an unstable and mixed form. Our work, in its form enable the understanding
 561 of different learning regimes by exploring initialization properties beyond just *absolute scale* and
 562 demonstrate that this transition can be accessed and controlled by adjusting a key parameter: the
 563 *relative scale*. Other studies, such as Kunin et al. (2024) and Xu & Zheng (2024), have similarly
 564 relaxed the balancedness assumption but were limited to single-output neuron settings.

565 C Appendix: Exact learning dynamics with prior knowledge

566 C.1 Appendix: Fukumizu Approach

567 **Lemma C.1.** *We introduce the variables*

$$\mathbf{Q} = \begin{bmatrix} \mathbf{W}_1^T \\ \mathbf{W}_2 \end{bmatrix} \quad \text{and} \quad \mathbf{Q} \mathbf{Q}^T = \begin{bmatrix} \mathbf{W}_1^T \mathbf{W}_1 & \mathbf{W}_1^T \mathbf{W}_2^T \\ \mathbf{W}_2 \mathbf{W}_1 & \mathbf{W}_2 \mathbf{W}_2^T \end{bmatrix}. \quad (16)$$

568 *Defining*

$$\mathbf{F} = \begin{bmatrix} -\frac{\lambda}{2} I & (\tilde{\Sigma}^{yx})^T \\ \tilde{\Sigma}^{yx} & \frac{\lambda}{2} I \end{bmatrix}, \quad (17)$$

569 *the gradient flow dynamics of $\mathbf{Q} \mathbf{Q}^T(t)$ can be written as a differential matrix Riccati equation*

$$\tau \frac{d}{dt} (\mathbf{Q} \mathbf{Q}^T) = \mathbf{F} \mathbf{Q} \mathbf{Q}^T + \mathbf{Q} \mathbf{Q}^T \mathbf{F} - (\mathbf{Q} \mathbf{Q}^T)^2. \quad (18)$$

570 *Proof.* We introduce the variables

$$\mathbf{Q} = \begin{bmatrix} \mathbf{W}_1^T \\ \mathbf{W}_2 \end{bmatrix} \quad \text{and} \quad \mathbf{Q} \mathbf{Q}^T = \begin{bmatrix} \mathbf{W}_1^T \mathbf{W}_1 & \mathbf{W}_1^T \mathbf{W}_2^T \\ \mathbf{W}_2 \mathbf{W}_1 & \mathbf{W}_2 \mathbf{W}_2^T \end{bmatrix}. \quad (19)$$

571 We compute the time derivative

$$\tau \frac{d}{dt} (\mathbf{Q} \mathbf{Q}^T) = \tau \begin{bmatrix} \frac{d\mathbf{W}_1^T}{dt} \mathbf{W}_1 + \mathbf{W}_1^T \frac{d\mathbf{W}_1}{dt} & \frac{d\mathbf{W}_1^T}{dt} \mathbf{W}_2 + \mathbf{W}_1^T \frac{d\mathbf{W}_2}{dt} \\ \frac{d\mathbf{W}_2}{dt} \mathbf{W}_1 + \mathbf{W}_2 \frac{d\mathbf{W}_1}{dt} & \frac{d\mathbf{W}_2}{dt} \mathbf{W}_2 + \mathbf{W}_2 \frac{d\mathbf{W}_2}{dt} \end{bmatrix}. \quad (20)$$

572 Using equations 18 and 19, we compute the four quadrants separately giving

$$\tau \left(\frac{d\mathbf{W}_1^T}{dt} \mathbf{W}_1 + \mathbf{W}_1^T \frac{d\mathbf{W}_1}{dt} \right) = \quad (21)$$

573

$$= (\Sigma^{yx} - \mathbf{W}_2 \mathbf{W}_1)^T \mathbf{W}_2 \mathbf{W}_1 + \mathbf{W}_1^T \mathbf{W}_2^T (\Sigma^{yx} - \mathbf{W}_2 \mathbf{W}_1) \quad (22)$$

574

$$= (\Sigma^{yx})^T \mathbf{W}_2 \mathbf{W}_1 + \mathbf{W}_1^T \mathbf{W}_2^T \Sigma^{yx} - \mathbf{W}_1^T \mathbf{W}_2^T \mathbf{W}_2 \mathbf{W}_1 - (\mathbf{W}_2 \mathbf{W}_1)^T \mathbf{W}_2 \mathbf{W}_1 \quad (23)$$

575

$$= (\Sigma^{yx})^T \mathbf{W}_2 \mathbf{W}_1 + \mathbf{W}_1^T \mathbf{W}_2^T \Sigma^{yx} - \mathbf{W}_1^T \mathbf{W}_2^T \mathbf{W}_2 \mathbf{W}_1 - \mathbf{W}_1^T \mathbf{W}_1 \mathbf{W}_1^T \mathbf{W}_1 - \lambda \mathbf{W}_1^T \mathbf{W}_1, \quad (24)$$

$$\tau \left(\frac{d\mathbf{W}_1^T}{dt} \mathbf{W}_2^T + \mathbf{W}_1^T \frac{d\mathbf{W}_2^T}{dt} \right) = \quad (25)$$

576

$$= (\Sigma^{yx} - \mathbf{W}_2 \mathbf{W}_1)^T \mathbf{W}_2 \mathbf{W}_2^T + \mathbf{W}_1^T \mathbf{W}_1 (\Sigma^{yx} - \mathbf{W}_2 \mathbf{W}_1)^T \quad (26)$$

577

$$= (\Sigma^{yx})^T \mathbf{W}_2 \mathbf{W}_2^T + \mathbf{W}_1^T \mathbf{W}_1 (\Sigma^{yx})^T - \mathbf{W}_1^T \mathbf{W}_1 \mathbf{W}_1^T \mathbf{W}_2^T - \mathbf{W}_1^T \mathbf{W}_2^T \mathbf{W}_2 \mathbf{W}_2^T, \quad (27)$$

$$\tau \left(\frac{d\mathbf{W}_2}{dt} \mathbf{W}_1 + \mathbf{W}_2 \frac{d\mathbf{W}_1}{dt} \right) = \quad (28)$$

578

$$= (\Sigma^{yx} - \mathbf{W}_2 \mathbf{W}_1) \mathbf{W}_1^T \mathbf{W}_1 + \mathbf{W}_2 \mathbf{W}_2^T (\Sigma^{yx} - \mathbf{W}_2 \mathbf{W}_1) \quad (29)$$

579

$$= \Sigma^{yx} \mathbf{W}_1^T \mathbf{W}_1 + \mathbf{W}_2 \mathbf{W}_2^T \Sigma^{yx} - \mathbf{W}_2 \mathbf{W}_2^T \mathbf{W}_2 \mathbf{W}_1 - \mathbf{W}_2 \mathbf{W}_1 \mathbf{W}_1^T \mathbf{W}_1, \quad (30)$$

$$\tau \left(\frac{d\mathbf{W}_2}{dt} \mathbf{W}_2^T + \mathbf{W}_2 \frac{d\mathbf{W}_2^T}{dt} \right) = \quad (31)$$

580

$$(\tilde{\Sigma}^{yx} - \mathbf{W}_2 \mathbf{W}_1) \mathbf{W}_1^T \mathbf{W}_2^T + \mathbf{W}_2 \mathbf{W}_1 (\tilde{\Sigma}^{yx} - \mathbf{W}_2 \mathbf{W}_1)^T \quad (32)$$

581

$$= \tilde{\Sigma}^{yx} \mathbf{W}_1^T \mathbf{W}_2^T + \mathbf{W}_2 \mathbf{W}_1 (\tilde{\Sigma}^{yx})^T - \mathbf{W}_2 \mathbf{W}_1 \mathbf{W}_1^T \mathbf{W}_2^T - \mathbf{W}_2 \mathbf{W}_1 (\mathbf{W}_2 \mathbf{W}_1)^T \quad (33)$$

582

$$= \tilde{\Sigma}^{yx} \mathbf{W}_1^T \mathbf{W}_2^T + \mathbf{W}_2 \mathbf{W}_1 (\tilde{\Sigma}^{yx})^T - \mathbf{W}_2 \mathbf{W}_1 \mathbf{W}_1^T \mathbf{W}_2^T - \mathbf{W}_2 \mathbf{W}_1 \mathbf{W}_1^T \mathbf{W}_2^T \quad (34)$$

583

$$= \tilde{\Sigma}^{yx} \mathbf{W}_1^T \mathbf{W}_2^T + \mathbf{W}_2 \mathbf{W}_1 (\tilde{\Sigma}^{yx})^T - \mathbf{W}_2 \mathbf{W}_1 \mathbf{W}_1^T \mathbf{W}_2^T - \mathbf{W}_2 \mathbf{W}_2^T \mathbf{W}_2 \mathbf{W}_2^T + \lambda \mathbf{W}_2 \mathbf{W}_2^T. \quad (35)$$

584 Defining

$$\mathbf{F} = \begin{bmatrix} -\frac{\lambda}{2} I & (\tilde{\Sigma}^{yx})^T \\ \tilde{\Sigma}^{yx} & \frac{\lambda}{2} I \end{bmatrix}, \quad (36)$$

585 the gradient flow dynamics of $\mathbf{Q}\mathbf{Q}^T(t)$ can be written as a differential matrix Riccati equation

$$\tau \frac{d}{dt} (\mathbf{Q}\mathbf{Q}^T) = \mathbf{F}\mathbf{Q}\mathbf{Q}^T + \mathbf{Q}\mathbf{Q}^T \mathbf{F} - (\mathbf{Q}\mathbf{Q}^T)^2. \quad (37)$$

586 We write $\tau \frac{d}{dt} (\mathbf{Q}\mathbf{Q}^T)$ for completeness

$$\begin{aligned} \tau \frac{d}{dt} (\mathbf{Q}\mathbf{Q}^T) &= \begin{bmatrix} -\frac{\lambda}{2} & (\tilde{\Sigma}^{yx})^T \\ \tilde{\Sigma}^{yx} & \frac{\lambda}{2} \end{bmatrix} \begin{bmatrix} \mathbf{W}_1^T \mathbf{W}_1 & \mathbf{W}_1^T \mathbf{W}_2 \\ \mathbf{W}_2 \mathbf{W}_1 & \mathbf{W}_2 \mathbf{W}_2^T \end{bmatrix} + \begin{bmatrix} \mathbf{W}_1^T \mathbf{W}_1 & \mathbf{W}_1^T \mathbf{W}_2 \\ \mathbf{W}_2 \mathbf{W}_1 & \mathbf{W}_2 \mathbf{W}_2^T \end{bmatrix}^T \begin{bmatrix} -\frac{\lambda}{2} & (\tilde{\Sigma}^{yx})^T \\ \tilde{\Sigma}^{yx} & \frac{\lambda}{2} \end{bmatrix} \\ &\quad - \begin{bmatrix} \mathbf{W}_1^T \mathbf{W}_1 & \mathbf{W}_1^T \mathbf{W}_2 \\ \mathbf{W}_2 \mathbf{W}_1 & \mathbf{W}_2 \mathbf{W}_2^T \end{bmatrix}^2 \end{aligned} \quad (38)$$

$$= \begin{bmatrix} -\frac{\lambda}{2} & (\tilde{\Sigma}^{yx})^T \\ \tilde{\Sigma}^{yx} & \frac{\lambda}{2} \end{bmatrix} \begin{bmatrix} \mathbf{W}_1^T \mathbf{W}_1 & \mathbf{W}_1^T \mathbf{W}_2 \\ \mathbf{W}_2 \mathbf{W}_1 & \mathbf{W}_2 \mathbf{W}_2^T \end{bmatrix} + \begin{bmatrix} \mathbf{W}_1^T \mathbf{W}_1 & \mathbf{W}_1^T \mathbf{W}_2 \\ \mathbf{W}_2 \mathbf{W}_1 & \mathbf{W}_2 \mathbf{W}_2^T \end{bmatrix}^T \begin{bmatrix} -\frac{\lambda}{2} & (\tilde{\Sigma}^{yx})^T \\ \tilde{\Sigma}^{yx} & \frac{\lambda}{2} \end{bmatrix} \quad (39)$$

$$\begin{aligned} &- \begin{bmatrix} \mathbf{W}_1^T \mathbf{W}_1 & \mathbf{W}_1^T \mathbf{W}_2 \\ \mathbf{W}_2 \mathbf{W}_1 & \mathbf{W}_2 \mathbf{W}_2^T \end{bmatrix} \begin{bmatrix} \mathbf{W}_1^T \mathbf{W}_1 & \mathbf{W}_1^T \mathbf{W}_2 \\ \mathbf{W}_2 \mathbf{W}_1 & \mathbf{W}_2 \mathbf{W}_2^T \end{bmatrix} \\ &= \begin{bmatrix} -\frac{\lambda}{2} \mathbf{W}_1^T \mathbf{W}_1 + (\tilde{\Sigma}^{yx})^T \mathbf{W}_2 \mathbf{W}_1 & -\frac{\lambda}{2} \mathbf{W}_1^T \mathbf{W}_2 + (\tilde{\Sigma}^{yx})^T \mathbf{W}_2 \mathbf{W}_2^T \\ \tilde{\Sigma}^{yx} \mathbf{W}_1^T \mathbf{W}_1 + \frac{\lambda}{2} \mathbf{W}_2 \mathbf{W}_1 & \tilde{\Sigma}^{yx} \mathbf{W}_1^T \mathbf{W}_2 + \frac{\lambda}{2} \mathbf{W}_2 \mathbf{W}_2^T \end{bmatrix} \\ &+ \begin{bmatrix} -\frac{\lambda}{2} \mathbf{W}_1^T \mathbf{W}_1 + \mathbf{W}_1^T \mathbf{W}_1 (\tilde{\Sigma}^{yx})^T & \frac{\lambda}{2} \mathbf{W}_1^T \mathbf{W}_2 + \mathbf{W}_1^T \mathbf{W}_2 (\tilde{\Sigma}^{yx})^T \\ -\frac{\lambda}{2} \mathbf{W}_2^T \mathbf{W}_1 + \mathbf{W}_2 \mathbf{W}_1 (\tilde{\Sigma}^{yx})^T & \frac{\lambda}{2} \mathbf{W}_2^T \mathbf{W}_2 + \mathbf{W}_2 \mathbf{W}_2^T (\tilde{\Sigma}^{yx})^T \end{bmatrix} \quad (40) \\ &- \begin{bmatrix} \mathbf{W}_1^T \mathbf{W}_1 & \mathbf{W}_1^T \mathbf{W}_2 \\ \mathbf{W}_2 \mathbf{W}_1 & \mathbf{W}_2 \mathbf{W}_2^T \end{bmatrix} \begin{bmatrix} \mathbf{W}_1^T \mathbf{W}_1 & \mathbf{W}_1^T \mathbf{W}_2 \\ \mathbf{W}_2 \mathbf{W}_1 & \mathbf{W}_2 \mathbf{W}_2^T \end{bmatrix} \end{aligned}$$

$$\begin{aligned}
&= \begin{bmatrix} -\frac{\lambda}{2} \mathbf{W}_1^T \mathbf{W}_1 + (\tilde{\Sigma}^{yx})^T \mathbf{W}_2 \mathbf{W}_1 & -\frac{\lambda}{2} \mathbf{W}_1^T \mathbf{W}_2 + (\tilde{\Sigma}^{yx})^T \mathbf{W}_2 \mathbf{W}_2^T \\ \tilde{\Sigma}^{yx} \mathbf{W}_1^T \mathbf{W}_1 + \frac{\lambda}{2} \mathbf{W}_2 \mathbf{W}_1 & \tilde{\Sigma}^{yx} \mathbf{W}_1^T \mathbf{W}_2^T + \frac{\lambda}{2} \mathbf{W}_2 \mathbf{W}_2^T \end{bmatrix} \\
&+ \begin{bmatrix} -\frac{\lambda}{2} \mathbf{W}_1^T \mathbf{W}_1 + \mathbf{W}_1^T \mathbf{W}_1 (\tilde{\Sigma}^{yx})^T & \frac{\lambda}{2} \mathbf{W}_1^T \mathbf{W}_2 + \mathbf{W}_1^T \mathbf{W}_2 (\tilde{\Sigma}^{yx})^T \\ -\frac{\lambda}{2} \mathbf{W}_2^T \mathbf{W}_1 + \mathbf{W}_2 \mathbf{W}_1 (\tilde{\Sigma}^{yx})^T & \frac{\lambda}{2} \mathbf{W}_2^T \mathbf{W}_2 + \mathbf{W}_2 \mathbf{W}_2^T (\tilde{\Sigma}^{yx})^T \end{bmatrix} \\
&- \begin{bmatrix} \mathbf{W}_1^T \mathbf{W}_1 \mathbf{W}_1^T \mathbf{W}_1 + \mathbf{W}_1^T \mathbf{W}_2 \mathbf{W}_2^T \mathbf{W}_1 & \mathbf{W}_1^T \mathbf{W}_1 \mathbf{W}_1^T \mathbf{W}_2 + \mathbf{W}_1^T \mathbf{W}_2 \mathbf{W}_2^T \mathbf{W}_2 \\ \mathbf{W}_2 \mathbf{W}_1 \mathbf{W}_1^T \mathbf{W}_1 + \mathbf{W}_2 \mathbf{W}_2^T \mathbf{W}_2 \mathbf{W}_1 & \mathbf{W}_2 \mathbf{W}_1 \mathbf{W}_1^T \mathbf{W}_2 + \mathbf{W}_2 \mathbf{W}_2^T \mathbf{W}_2 \mathbf{W}_2^T \end{bmatrix}
\end{aligned} \tag{45}$$

587

□

588 The four quadrants of 20 are equivalent to equations 24, 27, 30, and 35 respectively.

589 C.2 $\mathbf{Q}\mathbf{Q}^T$ Diagonalisation

590 **Lemma C.2.** *If $\mathbf{F} = \mathbf{P}\mathbf{\Lambda}\mathbf{P}^T$ is symmetric and diagonalizable, then the matrix Riccati differential*
591 *equation $\frac{d}{dt}(\mathbf{Q}\mathbf{Q}^T) = \mathbf{F}\mathbf{Q}\mathbf{Q}^T + \mathbf{Q}\mathbf{Q}^T\mathbf{F} - (\mathbf{Q}\mathbf{Q}^T)^2$ with initialization $\mathbf{Q}\mathbf{Q}^T(0) = \mathbf{Q}(0)\mathbf{Q}(0)^T$*
592 *has a unique solution for all $t \geq 0$, and the solution is given by*

$$\mathbf{Q}\mathbf{Q}^T(t) = e^{\mathbf{F}\frac{t}{\tau}}\mathbf{Q}(0) \left[\mathbf{I} + \mathbf{Q}(0)^T \mathbf{P} \left(\frac{e^{2\mathbf{\Lambda}\frac{t}{\tau}} - \mathbf{I}}{2\mathbf{\Lambda}} \right) \mathbf{P}^T \mathbf{Q}(0) \right]^{-1} \mathbf{Q}(0)^T e^{\mathbf{F}\frac{t}{\tau}}. \tag{41}$$

593 *This is true even when there exists $\mathbf{\Lambda}_i = 0$.*

594 *Proof.* First we show that there exists a unique solution to the initial value problem stated. This is
595 true by Picard-Lindelöf theorem. Now we show that the provided solution satisfies the ODE. Let

596 $\mathbf{L} = e^{\mathbf{F}\frac{t}{\tau}}\mathbf{Q}(0)$ and $\mathbf{C} = \mathbf{I} + \mathbf{Q}(0)^T \mathbf{P} \left(\frac{e^{2\mathbf{\Lambda}\frac{t}{\tau}} - \mathbf{I}}{2\mathbf{\Lambda}} \right) \mathbf{P}^T \mathbf{Q}(0)$ such that solution $\mathbf{Q}\mathbf{Q}^T(t) = \mathbf{L}\mathbf{C}^{-1}\mathbf{L}^T$.

597 The time derivative of $\mathbf{Q}\mathbf{Q}^T$ is then given by

$$\frac{d}{dt}(\mathbf{Q}\mathbf{Q}^T) = \frac{d}{dt}(\mathbf{L})\mathbf{C}^{-1}\mathbf{L}^T + \mathbf{L} \frac{d}{dt}(\mathbf{C}^{-1})\mathbf{L}^T + \mathbf{L}\mathbf{C}^{-1} \frac{d}{dt}(\mathbf{L}^T) \tag{42}$$

598 Solving for these derivatives individually, we find

$$\frac{d}{dt}(\mathbf{L}) = \frac{d}{dt}e^{\mathbf{F}\frac{t}{\tau}}\mathbf{Q}(0) = \mathbf{F}e^{\mathbf{F}\frac{t}{\tau}}\mathbf{Q}(0) = \mathbf{F}\mathbf{L} \tag{43}$$

$$\frac{d}{dt}(\mathbf{C}^{-1}) = -\mathbf{C}^{-1} \frac{d}{dt}(\mathbf{C})\mathbf{C}^{-1} = -\mathbf{C}^{-1}\mathbf{Q}(0)^T \mathbf{P} \frac{d}{dt} \left(\frac{e^{2\mathbf{\Lambda}\frac{t}{\tau}} - \mathbf{I}}{2\mathbf{\Lambda}} \right) \mathbf{P}^T \mathbf{Q}(0)\mathbf{C}^{-1} \tag{44}$$

599 We consider the derivative of the fraction separately,

$$\frac{d}{dt} \left(\frac{e^{2\mathbf{\Lambda}\frac{t}{\tau}} - \mathbf{I}}{2\mathbf{\Lambda}} \right) = e^{2\mathbf{\Lambda}\frac{t}{\tau}} \tag{45}$$

600 this is true even in the limit as $\lambda_i \rightarrow 0$. Plugging these derivatives back in we see that the solution
601 satisfies the ODE. Lastly, let $t = 0$, we see that the the solution satisfies the initial conditions. □

602 C.3 \mathbf{F} Diagonalization

603 **Lemma C.3.** *Under assumptions of full-rank 4, the eigendecomposition of $\mathbf{F} = \mathbf{P}\mathbf{\Lambda}\mathbf{P}^T$ where*

$$\mathbf{P} = \frac{1}{\sqrt{2}} \begin{pmatrix} \tilde{\mathbf{V}}(\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}}) & \tilde{\mathbf{V}}(\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}}) & \sqrt{2}\tilde{\mathbf{V}}_{\perp} \\ \tilde{\mathbf{U}}(\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}}) & -\tilde{\mathbf{U}}(\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}}) & \sqrt{2}\tilde{\mathbf{U}}_{\perp} \end{pmatrix}, \quad \mathbf{\Lambda} = \begin{pmatrix} \tilde{\mathbf{S}}_{\lambda} & 0 & 0 \\ 0 & -\tilde{\mathbf{S}}_{\lambda} & 0 \\ 0 & 0 & \lambda_{\perp} \end{pmatrix} \tag{46}$$

604 *and the matrices $\tilde{\mathbf{S}}_{\lambda}$, λ_{\perp} , $\tilde{\mathbf{H}}$, and $\tilde{\mathbf{G}}$ are the diagonal matrices defined as:*

$$\tilde{\mathbf{S}}_{\lambda} = \sqrt{\tilde{\mathbf{S}}^2 + \frac{\lambda^2}{4}\mathbf{I}}, \quad \lambda_{\perp} = \text{sgn}(N_o - N_i) \frac{\lambda}{2} \mathbf{I}, \quad \tilde{\mathbf{H}} = \text{sgn}(\lambda) \sqrt{\frac{\tilde{\mathbf{S}}_{\lambda} - \tilde{\mathbf{S}}}{\tilde{\mathbf{S}}_{\lambda} + \tilde{\mathbf{S}}}}, \quad \tilde{\mathbf{G}} = \frac{1}{\sqrt{\mathbf{I} + \tilde{\mathbf{H}}^2}}. \tag{47}$$

605 Beyond the invertibility of F , notice from the equation (Fukumizu solution) we need to understand
 606 the relationship between F and $Q(0)$. To do this the following lemma relates the structure between
 607 the SVD of the model with the SVD structure of the individual parameters.

608 *Proof.* We leave for the reader by computing

$$\mathbf{F} = \mathbf{P}\mathbf{A}\mathbf{P}^T \quad (48)$$

609

□

610 C.4 Solution Unequal-Input-Output

611 **Theorem C.4.** Under the assumptions of whitened inputs, 1, lambda-balanced weights 2, no bot-
 612 tleneck 3, and full rank 4, the temporal dynamics of $\mathbf{Q}\mathbf{Q}^T$ are

$$\mathbf{Q}\mathbf{Q}^T(t) = \begin{pmatrix} \mathbf{Z}_1\mathbf{A}^{-1}\mathbf{Z}_1^T & \mathbf{Z}_1\mathbf{A}^{-1}\mathbf{Z}_2^T \\ \mathbf{Z}_2\mathbf{A}^{-1}\mathbf{Z}_1^T & \mathbf{Z}_2\mathbf{A}^{-1}\mathbf{Z}_2^T \end{pmatrix},$$

613 where the variables $\mathbf{Z}_1 \in \mathbb{R}^{N_i \times N_h}$, $\mathbf{Z}_2 \in \mathbb{R}^{N_o \times N_h}$, and $\mathbf{A} \in \mathbb{R}^{N_h \times N_h}$ are defined as

$$\mathbf{Z}_1(t) = \frac{1}{2}\tilde{\mathbf{V}}(\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}})e^{\tilde{S}\lambda\frac{t}{\tau}}\mathbf{B}^T - \frac{1}{2}\tilde{\mathbf{V}}(\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}})e^{-\tilde{S}\lambda\frac{t}{\tau}}\mathbf{C}^T + \tilde{\mathbf{V}}_{\perp}e^{\lambda_{\perp}\frac{t}{\tau}}\tilde{\mathbf{V}}_{\perp}^T\mathbf{W}_1(0)^T \quad (49)$$

$$\mathbf{Z}_2(t) = \frac{1}{2}\tilde{\mathbf{U}}(\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}})e^{\tilde{S}\lambda\frac{t}{\tau}}\mathbf{B}^T + \frac{1}{2}\tilde{\mathbf{U}}(\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}})e^{-\tilde{S}\lambda\frac{t}{\tau}}\mathbf{C}^T + \tilde{\mathbf{U}}_{\perp}e^{\lambda_{\perp}\frac{t}{\tau}}\tilde{\mathbf{U}}_{\perp}^T\mathbf{W}_2(0) \quad (50)$$

$$\begin{aligned} \mathbf{A}(t) = & \mathbf{I} + \mathbf{B} \left(\frac{e^{2\tilde{S}\lambda\frac{t}{\tau}} - \mathbf{I}}{4\tilde{S}\lambda} \right) \mathbf{B}^T - \mathbf{C} \left(\frac{e^{-2\tilde{S}\lambda\frac{t}{\tau}} - \mathbf{I}}{4\tilde{S}\lambda} \right) \mathbf{C}^T + \mathbf{W}_2(0)^T \tilde{\mathbf{U}}_{\perp} \left(\frac{e^{\lambda_{\perp}\frac{t}{\tau}} - \mathbf{I}}{\lambda_{\perp}} \right) \tilde{\mathbf{U}}_{\perp}^T \mathbf{W}_2(0) \\ & + \mathbf{W}_1(0) \tilde{\mathbf{V}}_{\perp} \left(\frac{e^{\lambda_{\perp}\frac{t}{\tau}} - \mathbf{I}}{\lambda_{\perp}} \right) \tilde{\mathbf{V}}_{\perp}^T \mathbf{W}_1(0)^T \end{aligned} \quad (51)$$

614 *Proof.* We start and use the diagonalization of \mathbf{F} to rewrite the matrix exponential of \mathbf{F} and \mathbf{F} . Note
 615 that $\mathbf{P}^T\mathbf{P} = \mathbf{P}\mathbf{P}^T = \mathbf{I}$ and therefore $\mathbf{P}^T = \mathbf{P}^{-1}$.

$$e^{\mathbf{F}\frac{t}{\tau}} = \mathbf{P}e^{\mathbf{\Gamma}}\mathbf{P}^T$$

$$\begin{aligned} &= \frac{1}{\sqrt{2}} \begin{bmatrix} \tilde{\mathbf{V}}(\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}}) & \tilde{\mathbf{V}}(\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}}) & \sqrt{2}\mathbf{V}_{\perp} \\ \tilde{\mathbf{U}}(\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}}) & -\tilde{\mathbf{U}}(\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}}) & \sqrt{2}\mathbf{U}_{\perp} \end{bmatrix} \begin{bmatrix} e^{\tilde{S}\lambda\frac{t}{\tau}} & 0 & 0 \\ 0 & e^{-\tilde{S}\lambda\frac{t}{\tau}} & 0 \\ 0 & 0 & e^{\lambda_{\perp}\frac{t}{\tau}} \end{bmatrix} \frac{1}{\sqrt{2}} \begin{bmatrix} \tilde{\mathbf{V}}(\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}}) & \tilde{\mathbf{V}}(\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}}) & \sqrt{2}\mathbf{V}_{\perp} \\ \tilde{\mathbf{U}}(\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}}) & -\tilde{\mathbf{U}}(\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}}) & \sqrt{2}\mathbf{U}_{\perp} \end{bmatrix}^T \\ &= \frac{1}{\sqrt{2}} \begin{bmatrix} \tilde{\mathbf{V}}(\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}}) & \tilde{\mathbf{V}}(\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}}) \\ \tilde{\mathbf{U}}(\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}}) & -\tilde{\mathbf{U}}(\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}}) \end{bmatrix} \begin{bmatrix} e^{\tilde{S}\lambda\frac{t}{\tau}} & 0 \\ 0 & e^{-\tilde{S}\lambda\frac{t}{\tau}} \end{bmatrix} \frac{1}{\sqrt{2}} \begin{bmatrix} \tilde{\mathbf{V}}(\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}}) & \tilde{\mathbf{V}}(\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}}) \\ \tilde{\mathbf{U}}(\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}}) & -\tilde{\mathbf{U}}(\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}}) \end{bmatrix}^T + 2\frac{1}{\sqrt{2}} \begin{bmatrix} \tilde{\mathbf{V}}_{\perp} \\ \tilde{\mathbf{U}}_{\perp} \end{bmatrix} e^{\lambda_{\perp}\frac{t}{\tau}} \frac{1}{\sqrt{2}} \begin{bmatrix} \tilde{\mathbf{V}}_{\perp} \\ \tilde{\mathbf{U}}_{\perp} \end{bmatrix}^T \\ &= \mathbf{O}e^{\Lambda\frac{t}{\tau}}\mathbf{O} + 2\mathbf{M}e^{\lambda_{\perp}\frac{t}{\tau}}\mathbf{M}^T. \end{aligned} \quad (52)$$

616

$$e^{\mathbf{F}\frac{t}{\tau}}\mathbf{F}^{-1}e^{\mathbf{F}\frac{t}{\tau}} - \mathbf{F}^{-1} = \mathbf{O}e^{\Lambda\frac{t}{\tau}}\mathbf{O}^T\mathbf{O}\mathbf{\Lambda}^{-1}\mathbf{O}^T\mathbf{O}e^{\Lambda\frac{t}{\tau}}\mathbf{O}^T - \mathbf{O}\mathbf{\Lambda}^{-1}\mathbf{O}^T + \mathbf{M}(e^{\lambda_{\perp}\frac{t}{\tau}} - \mathbf{I})(\lambda_{\perp})^{-1}\mathbf{M}^T. \quad (53)$$

$$\mathbf{F} = \mathbf{O}\mathbf{\Lambda}\mathbf{O}^T + 2\mathbf{M}\lambda_{\perp}\mathbf{M}^T \quad (54)$$

617 Where $\mathbf{M} = \frac{1}{\sqrt{2}} \begin{bmatrix} \tilde{\mathbf{V}}_{\perp} \\ \tilde{\mathbf{U}}_{\perp} \end{bmatrix}^T$. Placing these expressions into equation 41 gives

$$\begin{aligned} \mathbf{Q}\mathbf{Q}^T(t) = & \left[\mathbf{O}e^{\Lambda\frac{t}{\tau}}\mathbf{O}^T + 2\mathbf{M}e^{\lambda_{\perp}\frac{t}{\tau}}\mathbf{M}^T \right] \mathbf{Q}(0) \\ & \left[\mathbf{I} + \frac{1}{2}\mathbf{Q}(0)^T \left(\mathbf{O} \left(e^{2\Lambda\frac{t}{\tau}} - \mathbf{I} \right) \mathbf{\Lambda}^{-1}\mathbf{O}^T + \mathbf{M}(e^{\lambda_{\perp}\frac{t}{\tau}} - \mathbf{I})\lambda_{\perp}^{-1}\mathbf{M}^T \right) \mathbf{Q}(0) \right]^{-1} \\ & \mathbf{Q}(0)^T \left[\mathbf{O}e^{\Lambda\frac{t}{\tau}}\mathbf{O}^T + 2\mathbf{M}e^{\lambda_{\perp}\frac{t}{\tau}}\mathbf{M}^T \right]^T \end{aligned} \quad (55)$$

$$\mathbf{O}^T\mathbf{Q}(0) = \frac{1}{\sqrt{2}} \begin{pmatrix} \tilde{\mathbf{V}}(\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}}) & \tilde{\mathbf{V}}(\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}}) \\ \tilde{\mathbf{U}}(\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}}) & -\tilde{\mathbf{U}}(\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}}) \end{pmatrix}^T \begin{pmatrix} \mathbf{W}_1^T(0) \\ \mathbf{W}_2^T(0) \end{pmatrix}$$

$$\begin{aligned}
&= \frac{1}{\sqrt{2}} \left(\begin{array}{c} (\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}})\tilde{\mathbf{V}}^T\mathbf{W}_1^T(0) + (\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}})\tilde{\mathbf{U}}^T\mathbf{W}_2(0) \\ (\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}})\tilde{\mathbf{V}}^T\mathbf{W}_1^T(0) - (\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}})\tilde{\mathbf{U}}^T\mathbf{W}_2(0) \end{array} \right) \\
&= \frac{1}{\sqrt{2}} \begin{pmatrix} \mathbf{B}^T \\ -\mathbf{C}^T \end{pmatrix} \tag{56}
\end{aligned}$$

618 where

$$\mathbf{B} = \mathbf{W}_2(0)^T\tilde{\mathbf{U}}(\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}}) + \mathbf{W}_1(0)\tilde{\mathbf{V}}(\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}}) \in \mathbb{R}^{N_h \times N_h} \tag{57}$$

$$\mathbf{C} = \mathbf{W}_2(0)^T\tilde{\mathbf{U}}(\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}}) - \mathbf{W}_1(0)\tilde{\mathbf{V}}(\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}}) \in \mathbb{R}^{N_h \times N_h} \tag{58}$$

$$\begin{aligned}
\mathbf{O}e^{\Lambda t/\tau} &= \frac{1}{\sqrt{2}} \begin{pmatrix} \tilde{\mathbf{V}}(\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}}) & \tilde{\mathbf{V}}(\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}}) \\ \tilde{\mathbf{U}}(\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}}) & -\tilde{\mathbf{U}}(\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}}) \end{pmatrix} \begin{pmatrix} e^{\tilde{\mathcal{S}}_\lambda \frac{t}{\tau}} & 0 \\ 0 & e^{-\tilde{\mathcal{S}}_\lambda \frac{t}{\tau}} \end{pmatrix} \\
&= \frac{1}{\sqrt{2}} \begin{pmatrix} \tilde{\mathbf{V}}(\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}})e^{\tilde{\mathcal{S}}_\lambda \frac{t}{\tau}} & \tilde{\mathbf{V}}(\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}})e^{-\tilde{\mathcal{S}}_\lambda \frac{t}{\tau}} \\ \tilde{\mathbf{U}}(\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}})e^{\tilde{\mathcal{S}}_\lambda \frac{t}{\tau}} & -\tilde{\mathbf{U}}(\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}})e^{-\tilde{\mathcal{S}}_\lambda \frac{t}{\tau}} \end{pmatrix} \tag{59}
\end{aligned}$$

$$\begin{aligned}
\mathbf{O}e^{\Lambda t/\tau}\mathbf{O}^T\mathbf{Q}(0) &= \frac{1}{2} \begin{pmatrix} \tilde{\mathbf{V}}(\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}})e^{\tilde{\mathcal{S}}_\lambda \frac{t}{\tau}} & \tilde{\mathbf{V}}(\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}})e^{-\tilde{\mathcal{S}}_\lambda \frac{t}{\tau}} \\ \tilde{\mathbf{U}}(\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}})e^{\tilde{\mathcal{S}}_\lambda \frac{t}{\tau}} & -\tilde{\mathbf{U}}(\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}})e^{-\tilde{\mathcal{S}}_\lambda \frac{t}{\tau}} \end{pmatrix} \begin{pmatrix} \mathbf{B}^T \\ -\mathbf{C}^T \end{pmatrix} \\
&= \frac{1}{2} \begin{pmatrix} \tilde{\mathbf{V}}(\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}})e^{\tilde{\mathcal{S}}_\lambda \frac{t}{\tau}}\mathbf{B}^T - \tilde{\mathbf{V}}(\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}})e^{-\tilde{\mathcal{S}}_\lambda \frac{t}{\tau}}\mathbf{C}^T \\ \tilde{\mathbf{U}}(\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}})e^{\tilde{\mathcal{S}}_\lambda \frac{t}{\tau}}\mathbf{B}^T + \tilde{\mathbf{U}}(\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}})e^{-\tilde{\mathcal{S}}_\lambda \frac{t}{\tau}}\mathbf{C}^T \end{pmatrix} \tag{60}
\end{aligned}$$

$$\begin{aligned}
2\mathbf{M}e^{\lambda_\perp \frac{t}{\tau}}\mathbf{M}^T\mathbf{Q}(0) &= 2\frac{1}{\sqrt{2}} \begin{bmatrix} \tilde{\mathbf{V}}_\perp \\ \tilde{\mathbf{U}}_\perp \end{bmatrix} \begin{bmatrix} e^{\lambda_\perp \frac{t}{\tau}} & 0 \\ 0 & e^{\lambda_\perp \frac{t}{\tau}} \end{bmatrix} \frac{1}{\sqrt{2}} \begin{bmatrix} \tilde{\mathbf{V}}_\perp \\ \tilde{\mathbf{U}}_\perp \end{bmatrix}^T \begin{bmatrix} \mathbf{W}_1(0)^T \\ \mathbf{W}_2(0) \end{bmatrix} \\
&= \begin{bmatrix} \tilde{\mathbf{V}}_\perp e^{\lambda_\perp \frac{t}{\tau}} \tilde{\mathbf{V}}_\perp^T & 0 \\ 0 & \tilde{\mathbf{U}}_\perp e^{\lambda_\perp \frac{t}{\tau}} \tilde{\mathbf{U}}_\perp^T \end{bmatrix} \begin{bmatrix} \mathbf{W}_1(0)^T \\ \mathbf{W}_2(0) \end{bmatrix} \\
&= \begin{bmatrix} \tilde{\mathbf{V}}_\perp e^{\lambda_\perp \frac{t}{\tau}} \tilde{\mathbf{V}}_\perp^T \mathbf{W}_1(0)^T \\ \tilde{\mathbf{U}}_\perp e^{\lambda_\perp \frac{t}{\tau}} \tilde{\mathbf{U}}_\perp^T \mathbf{W}_2(0) \end{bmatrix} \tag{61}
\end{aligned}$$

619 Putting it together we get the expressions for $\mathbf{Z}_1(t)$ and $\mathbf{Z}_2(t)$

$$\begin{aligned}
&\left[\mathbf{O}e^{\Lambda \frac{t}{\tau}}\mathbf{O}^T + 2\mathbf{M}e^{\lambda_\perp \frac{t}{\tau}}\mathbf{M}^T \right] \mathbf{Q}(0) = \\
&= \frac{1}{2} \begin{pmatrix} \tilde{\mathbf{V}}(\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}})e^{\tilde{\mathcal{S}}_\lambda \frac{t}{\tau}}\mathbf{B}^T - \tilde{\mathbf{V}}(\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}})e^{-\tilde{\mathcal{S}}_\lambda \frac{t}{\tau}}\mathbf{C}^T \\ \tilde{\mathbf{U}}(\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}})e^{\tilde{\mathcal{S}}_\lambda \frac{t}{\tau}}\mathbf{B}^T + \tilde{\mathbf{U}}(\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}})e^{-\tilde{\mathcal{S}}_\lambda \frac{t}{\tau}}\mathbf{C}^T \end{pmatrix} + \begin{bmatrix} \tilde{\mathbf{V}}_\perp e^{\lambda_\perp \frac{t}{\tau}} \tilde{\mathbf{V}}_\perp^T \mathbf{W}_1(0)^T \\ \tilde{\mathbf{U}}_\perp e^{\lambda_\perp \frac{t}{\tau}} \tilde{\mathbf{U}}_\perp^T \mathbf{W}_2(0) \end{bmatrix} \tag{62}
\end{aligned}$$

$$\mathbf{Z}_1(t) = \frac{1}{2}\tilde{\mathbf{V}}(\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}})e^{\tilde{\mathcal{S}}_\lambda \frac{t}{\tau}}\mathbf{B}^T - \frac{1}{2}\tilde{\mathbf{V}}(\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}})e^{-\tilde{\mathcal{S}}_\lambda \frac{t}{\tau}}\mathbf{C}^T + \tilde{\mathbf{V}}_\perp e^{\lambda_\perp \frac{t}{\tau}} \tilde{\mathbf{V}}_\perp^T \mathbf{W}_1(0)^T \tag{63}$$

$$\mathbf{Z}_2(t) = \frac{1}{2}\tilde{\mathbf{U}}(\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}})e^{\tilde{\mathcal{S}}_\lambda \frac{t}{\tau}}\mathbf{B}^T + \frac{1}{2}\tilde{\mathbf{U}}(\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}})e^{-\tilde{\mathcal{S}}_\lambda \frac{t}{\tau}}\mathbf{C}^T + \tilde{\mathbf{U}}_\perp e^{\lambda_\perp \frac{t}{\tau}} \tilde{\mathbf{U}}_\perp^T \mathbf{W}_2(0) \tag{64}$$

620 We now compute the terms inside the inverse

$$\begin{aligned}
& \mathbf{Q}(0)^T \mathbf{M} (e^{\lambda_{\perp} \frac{t}{\tau}}) \lambda_{\perp}^{-1} \mathbf{M}^T \mathbf{Q}(0) \\
&= [\mathbf{W}_1(0) \quad \mathbf{W}_2(0)^T] \frac{1}{\sqrt{2}} \begin{bmatrix} \tilde{\mathbf{V}}_{\perp} \\ \tilde{\mathbf{U}}_{\perp} \end{bmatrix} \begin{bmatrix} e^{\lambda_{\perp} \frac{t}{\tau}} & 0 \\ 0 & e^{\lambda_{\perp} \frac{t}{\tau}} \end{bmatrix} \begin{bmatrix} \lambda_{\perp} & 0 \\ 0 & \lambda_{\perp} \end{bmatrix}^{-1} \frac{1}{\sqrt{2}} \begin{bmatrix} \tilde{\mathbf{V}}_{\perp} \\ \tilde{\mathbf{U}}_{\perp} \end{bmatrix}^T \begin{bmatrix} \mathbf{W}_1(0)^T \\ \mathbf{W}_2(0) \end{bmatrix} \\
&= [\mathbf{W}_1(0) \quad \mathbf{W}_2(0)^T] \begin{bmatrix} e^{\lambda_{\perp} \frac{t}{\tau}} \lambda_{\perp}^{-1} \tilde{\mathbf{V}}_{\perp} \tilde{\mathbf{V}}_{\perp}^T \mathbf{W}_1(0)^T \\ e^{\lambda_{\perp} \frac{t}{\tau}} \lambda_{\perp}^{-1} \tilde{\mathbf{U}}_{\perp} \tilde{\mathbf{U}}_{\perp}^T \mathbf{W}_2(0) \end{bmatrix} \\
&= \left[\left(\mathbf{W}_1(0) \tilde{\mathbf{V}}_{\perp} e^{\lambda_{\perp} \frac{t}{\tau}} \lambda_{\perp}^{-1} \tilde{\mathbf{V}}_{\perp}^T \mathbf{W}_1(0)^T + \mathbf{W}_2(0)^T \tilde{\mathbf{U}}_{\perp} e^{\lambda_{\perp} \frac{t}{\tau}} \lambda_{\perp}^{-1} \tilde{\mathbf{U}}_{\perp}^T \mathbf{W}_2(0) \right) \right] \quad (65)
\end{aligned}$$

$$\begin{aligned}
\mathbf{Q}(0)^T \mathbf{M} \lambda_{\perp}^{-1} \mathbf{M}^T \mathbf{Q}(0) &= 2 [\mathbf{W}_1(0) \quad \mathbf{W}_2(0)^T] \frac{1}{\sqrt{2}} \begin{bmatrix} \tilde{\mathbf{V}}_{\perp} \\ \tilde{\mathbf{U}}_{\perp} \end{bmatrix} \begin{bmatrix} \lambda_{\perp} & 0 \\ 0 & \lambda_{\perp} \end{bmatrix}^{-1} \frac{1}{\sqrt{2}} \begin{bmatrix} \tilde{\mathbf{V}}_{\perp} \\ \tilde{\mathbf{U}}_{\perp} \end{bmatrix}^T \begin{bmatrix} \mathbf{W}_1(0)^T \\ \mathbf{W}_2(0) \end{bmatrix} \\
&= [\mathbf{W}_1(0) \quad \mathbf{W}_2(0)^T] \begin{bmatrix} \tilde{\mathbf{V}}_{\perp} \\ \tilde{\mathbf{U}}_{\perp} \end{bmatrix} \begin{bmatrix} \lambda_{\perp}^{-1} \tilde{\mathbf{V}}_{\perp} \tilde{\mathbf{V}}_{\perp}^T \mathbf{W}_1(0)^T \\ \lambda_{\perp}^{-1} \tilde{\mathbf{U}}_{\perp} \tilde{\mathbf{U}}_{\perp}^T \mathbf{W}_2(0) \end{bmatrix} \\
&= [\mathbf{W}_1(0) \tilde{\mathbf{V}}_{\perp} \lambda_{\perp}^{-1} \tilde{\mathbf{V}}_{\perp}^T \mathbf{W}_1(0)^T + \mathbf{W}_2(0)^T \tilde{\mathbf{U}}_{\perp} \lambda_{\perp}^{-1} \tilde{\mathbf{U}}_{\perp}^T \mathbf{W}_2(0)] \quad (66)
\end{aligned}$$

621 Now

$$\begin{aligned}
\frac{1}{2} \mathbf{Q}(0)^T \mathbf{O} \left(e^{2\Lambda \frac{t}{\tau}} - \mathbf{I} \right) \Lambda^{-1} \mathbf{O}^T &= \frac{1}{4} [\mathbf{B} - \mathbf{C}] \left(e^{\Lambda \frac{t}{\tau}} - \mathbf{I} \right) \Lambda^{-1} \begin{bmatrix} \mathbf{B}^T \\ -\mathbf{C}^T \end{bmatrix} \\
&= \frac{1}{4} \left(\mathbf{B} \left(e^{2\tilde{S}_{\lambda} \frac{t}{\tau}} - \mathbf{I} \right) (\tilde{S}_{\lambda})^{-1} \mathbf{B}^T - \mathbf{C} \left(e^{-2\tilde{S}_{\lambda} \frac{t}{\tau}} - \mathbf{I} \right) (\tilde{S}_{\lambda})^{-1} \mathbf{C}^T \right) \quad (67)
\end{aligned}$$

622 Putting it all together

$$\begin{aligned}
\mathbf{A}(t) &= \mathbf{I} + \mathbf{B} \left(\frac{e^{2\tilde{S}_{\lambda} \frac{t}{\tau}} - \mathbf{I}}{4\tilde{S}_{\lambda}} \right) \mathbf{B}^T - \mathbf{C} \left(\frac{e^{-2\tilde{S}_{\lambda} \frac{t}{\tau}} - \mathbf{I}}{4\tilde{S}_{\lambda}} \right) \mathbf{C}^T + \mathbf{W}_2(0)^T \tilde{\mathbf{U}}_{\perp} \left(\frac{e^{\lambda_{\perp} \frac{t}{\tau}} - \mathbf{I}}{\lambda_{\perp}} \right) \tilde{\mathbf{U}}_{\perp}^T \mathbf{W}_2(0) \\
&\quad + \mathbf{W}_1(0) \tilde{\mathbf{V}}_{\perp} \left(\frac{e^{\lambda_{\perp} \frac{t}{\tau}} - \mathbf{I}}{\lambda_{\perp}} \right) \tilde{\mathbf{V}}_{\perp}^T \mathbf{W}_1(0)^T \quad (68)
\end{aligned}$$

623 So, final form:

$$\begin{aligned}
& \mathbf{Q}\mathbf{Q}^T(t) = \\
& \left[\begin{aligned} & \left(\frac{1}{2} \tilde{\mathbf{V}} (\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}}) e^{\tilde{S}_{\lambda} \frac{t}{\tau}} \mathbf{B}^T - \frac{1}{2} \tilde{\mathbf{V}} (\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}}) e^{-\tilde{S}_{\lambda} \frac{t}{\tau}} \mathbf{C}^T + \tilde{\mathbf{V}}_{\perp} e^{\lambda_{\perp} \frac{t}{\tau}} \tilde{\mathbf{V}}_{\perp}^T \mathbf{W}_1(0)^T \right) \\ & \left(\frac{1}{2} \tilde{\mathbf{U}} (\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}}) e^{\tilde{S}_{\lambda} \frac{t}{\tau}} \mathbf{B}^T + \frac{1}{2} \tilde{\mathbf{U}} (\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}}) e^{-\tilde{S}_{\lambda} \frac{t}{\tau}} \mathbf{C}^T + \tilde{\mathbf{U}}_{\perp} e^{\lambda_{\perp} \frac{t}{\tau}} \tilde{\mathbf{U}}_{\perp}^T \mathbf{W}_2(0) \right) \end{aligned} \right] \\
& \left[\mathbf{I} + \frac{1}{4} \left(\mathbf{B} \left(\frac{e^{2\tilde{S}_{\lambda} \frac{t}{\tau}} - \mathbf{I}}{\tilde{S}_{\lambda}} \right) \mathbf{B}^T - \mathbf{C} \left(\frac{e^{-2\tilde{S}_{\lambda} \frac{t}{\tau}} - \mathbf{I}}{\tilde{S}_{\lambda}} \right) \mathbf{C}^T \right) \right. \\
& \left. + \mathbf{W}_2(0)^T \tilde{\mathbf{U}}_{\perp} \left(\frac{e^{\lambda_{\perp} \frac{t}{\tau}} - \mathbf{I}}{\lambda_{\perp}} \right) \tilde{\mathbf{U}}_{\perp}^T \mathbf{W}_2(0) + \mathbf{W}_1(0) \tilde{\mathbf{V}}_{\perp} \left(\frac{e^{\lambda_{\perp} \frac{t}{\tau}} - \mathbf{I}}{\lambda_{\perp}} \right) \tilde{\mathbf{V}}_{\perp}^T \mathbf{W}_1(0)^T \right]^{-1} \\
& \left[\begin{aligned} & \left(\frac{1}{2} \tilde{\mathbf{V}} (\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}}) e^{\tilde{S}_{\lambda} \frac{t}{\tau}} \mathbf{B}^T - \frac{1}{2} \tilde{\mathbf{V}} (\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}}) e^{-\tilde{S}_{\lambda} \frac{t}{\tau}} \mathbf{C}^T + \tilde{\mathbf{V}}_{\perp} e^{\lambda_{\perp} \frac{t}{\tau}} \tilde{\mathbf{V}}_{\perp}^T \mathbf{W}_1(0)^T \right) \\ & \left(\frac{1}{2} \tilde{\mathbf{U}} (\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}}) e^{\tilde{S}_{\lambda} \frac{t}{\tau}} \mathbf{B}^T + \frac{1}{2} \tilde{\mathbf{U}} (\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}}) e^{-\tilde{S}_{\lambda} \frac{t}{\tau}} \mathbf{C}^T + \tilde{\mathbf{U}}_{\perp} e^{\lambda_{\perp} \frac{t}{\tau}} \tilde{\mathbf{U}}_{\perp}^T \mathbf{W}_2(0) \right) \end{aligned} \right]^T \quad (69)
\end{aligned}$$

624

□

625 C.5 Stable solution Unequal-Input-Output

626 **Theorem C.5.** Given the assumptions of Theorem 2.3 further assuming that \mathbf{B} is invertible and
627 defining $e^{\lambda_{\perp} \frac{t}{\tau}} = \text{sgn}(N_o - N_i) \frac{\lambda}{2}$, the temporal evolution of $\mathbf{Q}\mathbf{Q}^T$ is described as follows:

$$\begin{aligned}
\mathbf{Q}\mathbf{Q}^T(t) &= \mathbf{Z} \left[e^{-\tilde{S}_\lambda \frac{t}{\tau}} \mathbf{B}^{-1} \mathbf{B}^{-T} e^{-\tilde{S}_\lambda \frac{t}{\tau}} \right. \\
&+ \left(\frac{\mathbf{I} - e^{-2\tilde{S}_\lambda \frac{t}{\tau}}}{4\tilde{S}_\lambda} \right) - e^{-\tilde{S}_\lambda \frac{t}{\tau}} \mathbf{B}^{-1} \mathbf{C} \left(\frac{e^{-\tilde{S}_\lambda \frac{t}{\tau}} - \mathbf{I}}{4\tilde{S}_\lambda} \right) \mathbf{C}^T \mathbf{B}^{-T} e^{-\tilde{S}_\lambda \frac{t}{\tau}} \\
&- e^{-\tilde{S}_\lambda \frac{t}{\tau}} \mathbf{B}^{-1} \mathbf{W}_2(0)^T \tilde{\mathbf{U}}_\perp \lambda_\perp^{-1} \tilde{\mathbf{U}}_\perp^T \mathbf{W}_2(0) \mathbf{B}^{-T} e^{-\tilde{S}_\lambda \frac{t}{\tau}} \\
&e^{-\tilde{S}_\lambda \frac{t}{\tau}} e^{\frac{\lambda_\perp}{2} \frac{t}{\tau}} \mathbf{B}^{-1} \mathbf{W}_2(0)^T \tilde{\mathbf{U}}_\perp \lambda_\perp^{-1} \tilde{\mathbf{U}}_\perp^T \mathbf{W}_2(0) \mathbf{B}^{-T} e^{-\tilde{S}_\lambda \frac{t}{\tau}} \\
&+ e^{-\tilde{S}_\lambda \frac{t}{\tau}} e^{\frac{\lambda}{2} \frac{t}{\tau}} \mathbf{B}^{-1} \mathbf{W}_1(0)^T \tilde{\mathbf{V}}_\perp \lambda_\perp^{-1} \tilde{\mathbf{V}}_\perp^T \mathbf{W}_1(0)^T \mathbf{B}^{-T} e^{-\tilde{S}_\lambda \frac{t}{\tau}} \\
&\left. - e^{-\tilde{S}_\lambda \frac{t}{\tau}} \mathbf{B}^{-1} \mathbf{W}_1(0)^T \tilde{\mathbf{V}}_\perp \lambda_\perp^{-1} \tilde{\mathbf{V}}_\perp^T \mathbf{W}_1(0)^T \mathbf{B}^{-T} e^{-\tilde{S}_\lambda \frac{t}{\tau}} \right]^{-1} \mathbf{Z}^T
\end{aligned} \tag{70}$$

628

$$\mathbf{Z} = \left(\begin{array}{l} \frac{1}{2} \tilde{\mathbf{V}} \left[(\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}}) - (\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}}) e^{-\tilde{S}_\lambda \frac{t}{\tau}} \mathbf{C}^T \mathbf{B}^{-T} e^{-\tilde{S}_\lambda \frac{t}{\tau}} \right] + \tilde{\mathbf{V}}_\perp \tilde{\mathbf{V}}_\perp^T \mathbf{W}_1(0) \mathbf{B}^{-T} e^{\lambda_\perp \frac{t}{\tau}} e^{-\tilde{S}_\lambda \frac{t}{\tau}} \\ \frac{1}{2} \tilde{\mathbf{U}} \left[(\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}}) + (\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}}) e^{-\tilde{S}_\lambda \frac{t}{\tau}} \mathbf{C}^T \mathbf{B}^{-T} e^{-\tilde{S}_\lambda \frac{t}{\tau}} \right] + \tilde{\mathbf{U}}_\perp \tilde{\mathbf{U}}_\perp^T \mathbf{W}_2(0)^T \mathbf{B}^{-T} e^{\lambda_\perp \frac{t}{\tau}} e^{-\tilde{S}_\lambda \frac{t}{\tau}} \end{array} \right) \tag{71}$$

629 *Proof.* We start from

$$\begin{aligned}
\mathbf{Q}\mathbf{Q}^T(t) &= \\
&\left[\left(\begin{array}{l} \frac{1}{2} \tilde{\mathbf{V}} (\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}}) e^{\tilde{S}_\lambda \frac{t}{\tau}} \mathbf{B}^T - \frac{1}{2} \tilde{\mathbf{V}} (\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}}) e^{-\tilde{S}_\lambda \frac{t}{\tau}} \mathbf{C}^T + \tilde{\mathbf{V}}_\perp e^{\lambda_\perp \frac{t}{\tau}} \tilde{\mathbf{V}}_\perp^T \mathbf{W}_1(0)^T \\ \frac{1}{2} \tilde{\mathbf{U}} (\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}}) e^{\tilde{S}_\lambda \frac{t}{\tau}} \mathbf{B}^T + \frac{1}{2} \tilde{\mathbf{U}} (\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}}) e^{-\tilde{S}_\lambda \frac{t}{\tau}} \mathbf{C}^T + \tilde{\mathbf{U}}_\perp e^{\lambda_\perp \frac{t}{\tau}} \tilde{\mathbf{U}}_\perp^T \mathbf{W}_2(0) \end{array} \right) \right] \\
&\left[\mathbf{I} + \frac{1}{4} \left(\mathbf{B} \left(\frac{e^{2\tilde{S}_\lambda \frac{t}{\tau}} - \mathbf{I}}{\tilde{S}_\lambda} \right) \mathbf{B}^T - \mathbf{C} \left(\frac{e^{-2\tilde{S}_\lambda \frac{t}{\tau}} - \mathbf{I}}{\tilde{S}_\lambda} \right) \mathbf{C}^T \right) \right. \\
&+ \mathbf{W}_2(0)^T \tilde{\mathbf{U}}_\perp \left(\frac{e^{\lambda_\perp \frac{t}{\tau}} - \mathbf{I}}{\lambda_\perp} \right) \tilde{\mathbf{U}}_\perp^T \mathbf{W}_2(0) + \mathbf{W}_1(0)^T \tilde{\mathbf{V}}_\perp \left(\frac{e^{\lambda_\perp \frac{t}{\tau}} - \mathbf{I}}{\lambda_\perp} \right) \tilde{\mathbf{V}}_\perp^T \mathbf{W}_1(0)^T \left. \right]^{-1} \\
&\left[\left(\begin{array}{l} \frac{1}{2} \tilde{\mathbf{V}} (\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}}) e^{\tilde{S}_\lambda \frac{t}{\tau}} \mathbf{B}^T - \frac{1}{2} \tilde{\mathbf{V}} (\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}}) e^{-\tilde{S}_\lambda \frac{t}{\tau}} \mathbf{C}^T + \tilde{\mathbf{V}}_\perp e^{\lambda_\perp \frac{t}{\tau}} \tilde{\mathbf{V}}_\perp^T \mathbf{W}_1(0)^T \\ \frac{1}{2} \tilde{\mathbf{U}} (\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}}) e^{\tilde{S}_\lambda \frac{t}{\tau}} \mathbf{B}^T + \frac{1}{2} \tilde{\mathbf{U}} (\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}}) e^{-\tilde{S}_\lambda \frac{t}{\tau}} \mathbf{C}^T + \tilde{\mathbf{U}}_\perp e^{\lambda_\perp \frac{t}{\tau}} \tilde{\mathbf{U}}_\perp^T \mathbf{W}_2(0) \end{array} \right) \right]^T
\end{aligned} \tag{72}$$

630 We extract $\mathbf{B}^{-T} e^{-\tilde{S}_\lambda \frac{t}{\tau}}$ from all terms as exemplified bellow

$$\mathbf{O} e^{\Lambda t/\tau} \mathbf{O}^T \mathbf{Q}(0) = \frac{1}{2} \left(\begin{array}{l} \tilde{\mathbf{V}} \left[(\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}}) - (\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}}) e^{-\tilde{S}_\lambda \frac{t}{\tau}} \mathbf{C}^T \mathbf{B}^{-T} e^{-\tilde{S}_\lambda \frac{t}{\tau}} \right] \\ \tilde{\mathbf{U}} \left[(\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}}) + (\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}}) e^{-\tilde{S}_\lambda \frac{t}{\tau}} \mathbf{C}^T \mathbf{B}^{-T} e^{-\tilde{S}_\lambda \frac{t}{\tau}} \right] \end{array} \right) \mathbf{B}^T e^{\tilde{S}_\lambda \frac{t}{\tau}} \tag{73}$$

634 **C.5.1 Proof Exact learning dynamics with prior knowledge unequal dimension**

635 We follow a similar derivation presented in Braun et al. (2022) and start with the following equation

$$\begin{aligned}
\mathbf{Q}\mathbf{Q}^T(t) &= \underbrace{\left[\mathbf{O}e^{\Lambda \frac{t}{\tau}} \mathbf{O}^T + 2\mathbf{M}e^{\lambda_{\perp} \frac{t}{\tau}} \mathbf{M}^T \right]}_{\mathbf{L}} \mathbf{Q}(0) \\
&\quad \underbrace{\left[\mathbf{I} + \frac{1}{2} \mathbf{Q}(0)^T \left(\mathbf{O} \left(e^{2\Lambda \frac{t}{\tau}} - \mathbf{I} \right) \Lambda^{-1} \mathbf{O}^T + \mathbf{M} \left(e^{\lambda_{\perp} \frac{t}{\tau}} - \mathbf{I} \right) \lambda_{\perp}^{-1} \mathbf{M}^T \right) \mathbf{Q}(0) \right]^{-1}}_{\mathbf{C}^{-1}} \\
&\quad \underbrace{\mathbf{Q}(0)^T \left[\mathbf{O}e^{\Lambda \frac{t}{\tau}} \mathbf{O}^T + 2\mathbf{M}e^{\lambda_{\perp} \frac{t}{\tau}} \mathbf{M}^T \right]}_{\mathbf{R}} \\
&= \mathbf{L}\mathbf{C}^{-1}\mathbf{R}, \tag{76}
\end{aligned}$$

636 Substituting our solution into the matrix Riccati equation then yields

$$\tau \frac{d}{dt} \mathbf{Q}\mathbf{Q}^T = \mathbf{F}\mathbf{Q}\mathbf{Q}^T + \mathbf{Q}\mathbf{Q}^T\mathbf{F} - (\mathbf{Q}\mathbf{Q}^T)^2 \tag{78}$$

$$\Rightarrow \tau \frac{d}{dt} \mathbf{L}\mathbf{C}^{-1}\mathbf{R} \stackrel{?}{=} \mathbf{F}\mathbf{L}\mathbf{C}^{-1}\mathbf{R} + \mathbf{L}\mathbf{C}^{-1}\mathbf{R}\mathbf{F} - \mathbf{L}\mathbf{C}^{-1}\mathbf{R}\mathbf{L}\mathbf{C}^{-1}\mathbf{R}. \tag{79}$$

637 Using the chain rule $\partial(\mathbf{A}\mathbf{B}) = (\partial\mathbf{A})\mathbf{B} + \mathbf{A}(\partial\mathbf{B})$ and the identities

$$\frac{d}{dt}(\mathbf{A}^{-1}) = \mathbf{A}^{-1} \left(\frac{d}{dt} \mathbf{A} \right) \mathbf{A}^{-1} \quad \text{and} \quad \frac{d}{dt}(e^{t\mathbf{A}}) = \mathbf{A}e^{t\mathbf{A}} = e^{t\mathbf{A}}\mathbf{A} \tag{80}$$

$$\tau \frac{d}{dt} \mathbf{Q}\mathbf{Q}^T = \tau \frac{d}{dt} \mathbf{L}\mathbf{C}^{-1}\mathbf{R} \tag{81}$$

$$= \tau \left(\frac{d}{dt} \mathbf{L} \right) \mathbf{C}^{-1}\mathbf{R} + \tau \mathbf{L} \left(\frac{d}{dt} \mathbf{C}^{-1}\mathbf{R} \right) \tag{82}$$

$$= \tau \left(\frac{d}{dt} \mathbf{L} \right) \mathbf{C}^{-1}\mathbf{R} + \tau \mathbf{L}\mathbf{C}^{-1} \left(\frac{d}{dt} \mathbf{R} \right) + \tau \mathbf{L} \left(\frac{d}{dt} \mathbf{C}^{-1} \right) \mathbf{R}, \tag{83}$$

638 Next, we note that

$$\mathbf{O} = \frac{1}{\sqrt{2}} \begin{pmatrix} \tilde{\mathbf{V}}(\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}}) & \tilde{\mathbf{V}}(\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}}) \\ \tilde{\mathbf{U}}(\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}}) & -\tilde{\mathbf{U}}(\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}}) \end{pmatrix}^T \tag{84}$$

639

$$\mathbf{O}^T\mathbf{O} = \frac{1}{\sqrt{2}} \begin{pmatrix} \tilde{\mathbf{V}}(\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}}) & \tilde{\mathbf{V}}(\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}}) \\ \tilde{\mathbf{U}}(\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}}) & -\tilde{\mathbf{U}}(\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}}) \end{pmatrix}^T \frac{1}{\sqrt{2}} \begin{pmatrix} \tilde{\mathbf{V}}(\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}}) & \tilde{\mathbf{V}}(\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}}) \\ \tilde{\mathbf{U}}(\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}}) & -\tilde{\mathbf{U}}(\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}}) \end{pmatrix} \tag{85}$$

$$= \mathbf{I} \tag{86}$$

$$\mathbf{O}^T\mathbf{M} = \frac{1}{\sqrt{2}} \begin{bmatrix} \tilde{\mathbf{V}}(\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}}) & \tilde{\mathbf{V}}(\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}}) \\ \tilde{\mathbf{U}}(\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}}) & -\tilde{\mathbf{U}}(\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}}) \end{bmatrix} \frac{1}{\sqrt{2}} \begin{bmatrix} \tilde{\mathbf{V}}_{\perp} \\ \tilde{\mathbf{U}}_{\perp} \end{bmatrix} \tag{87}$$

$$= \frac{1}{2} \begin{bmatrix} (\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}})^T \tilde{\mathbf{V}}^T \tilde{\mathbf{V}}_{\perp} + (\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}})^T \tilde{\mathbf{U}}^T \tilde{\mathbf{U}}_{\perp} \\ (\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}})^T \tilde{\mathbf{V}}^T \tilde{\mathbf{V}}_{\perp} - (\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}})^T \tilde{\mathbf{U}}^T \tilde{\mathbf{U}}_{\perp} \end{bmatrix} \tag{88}$$

$$= \mathbf{0} \tag{89}$$

640 and

$$\mathbf{M}^T \mathbf{O} = \frac{1}{\sqrt{2}} \begin{bmatrix} \tilde{\mathbf{V}}_{\perp}^T & \tilde{\mathbf{U}}_{\perp}^T \end{bmatrix} \begin{pmatrix} \tilde{\mathbf{V}}(\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}}) & \tilde{\mathbf{V}}(\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}}) \\ \tilde{\mathbf{U}}(\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}}) & -\tilde{\mathbf{U}}(\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}}) \end{pmatrix} \quad (90)$$

$$= \frac{1}{2} \begin{bmatrix} \tilde{\mathbf{V}}_{\perp}^T \tilde{\mathbf{V}}(\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}}) + \tilde{\mathbf{U}}_{\perp}^T \tilde{\mathbf{U}}(\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}}) \\ \tilde{\mathbf{V}}_{\perp}^T \tilde{\mathbf{V}}(\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}}) - \tilde{\mathbf{U}}_{\perp}^T \tilde{\mathbf{U}}(\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}}) \end{bmatrix} \quad (91)$$

$$= \mathbf{0}. \quad (92)$$

641 we get

$$\tau \frac{d}{dt} \mathbf{Q} \mathbf{Q}^T = \tau \frac{d}{dt} (\mathbf{L} \mathbf{C}^{-1} \mathbf{R}) \quad (93)$$

$$= \tau \left(\frac{d}{dt} \mathbf{L} \right) \mathbf{C}^{-1} \mathbf{R} + \tau \mathbf{L} \left(\frac{d}{dt} \mathbf{C}^{-1} \mathbf{R} \right) \quad (94)$$

$$= \tau \left(\frac{d}{dt} \mathbf{L} \right) \mathbf{C}^{-1} \mathbf{R} + \tau \mathbf{L} \mathbf{C}^{-1} \left(\frac{d}{dt} \mathbf{R} \right) + \tau \mathbf{L} \left(\frac{d}{dt} \mathbf{C}^{-1} \right) \mathbf{R}, \quad (95)$$

642 with

$$\tau \left(\frac{d}{dt} \mathbf{L} \right) \mathbf{C}^{-1} \mathbf{R} = \tau \left(\mathbf{O} \frac{1}{\tau} \boldsymbol{\Lambda} e^{\boldsymbol{\Lambda} \frac{t}{\tau}} \mathbf{O}^T + 2\mathbf{M} \frac{\lambda_{\perp} \mathbf{I}}{2\tau} e^{\lambda_{\perp} \frac{t}{\tau}} \mathbf{M}^T \right) \mathbf{Q}(0) \mathbf{C}^{-1} \mathbf{R} \quad (96)$$

$$= \left(\mathbf{O} \boldsymbol{\Lambda} e^{\boldsymbol{\Lambda} \frac{t}{\tau}} \mathbf{O}^T + \mathbf{M} \lambda_{\perp} \mathbf{I} e^{\lambda_{\perp} \frac{t}{\tau}} \mathbf{M}^T \right) \mathbf{Q}(0) \mathbf{C}^{-1} \mathbf{R} \quad (97)$$

$$= \left(\mathbf{O} \lambda_{\perp} \mathbf{O}^T + 2\mathbf{M} \lambda_{\perp} \mathbf{M}^T \right) \left(\mathbf{O} e^{\boldsymbol{\Lambda} \frac{t}{\tau}} \mathbf{O}^T + 2\mathbf{M} e^{\lambda_{\perp} \frac{t}{\tau}} \mathbf{M}^T \right) \mathbf{Q}(0) \mathbf{C}^{-1} \mathbf{R} \quad (98)$$

$$= \mathbf{F} \mathbf{L} \mathbf{C}^{-1} \mathbf{R}, \quad (99)$$

$$\tau \mathbf{L} \mathbf{C}^{-1} \left(\frac{d}{dt} \mathbf{R} \right) = \tau \mathbf{L} \mathbf{C}^{-1} \mathbf{Q}(0)^T \left(\mathbf{O} \frac{1}{\tau} e^{\boldsymbol{\Lambda} \frac{t}{\tau}} \boldsymbol{\Lambda} \mathbf{O}^T + 2\mathbf{M} e^{\lambda_{\perp} \frac{t}{\tau}} \frac{\lambda_{\perp} \mathbf{I}}{2\tau} \mathbf{M}^T \right) \quad (100)$$

$$= \mathbf{L} \mathbf{C}^{-1} \mathbf{Q}(0)^T \left(\mathbf{O} \frac{1}{\tau} e^{\boldsymbol{\Lambda} \frac{t}{\tau}} \boldsymbol{\Lambda} \mathbf{O}^T + 2\mathbf{M} e^{\lambda_{\perp} \frac{t}{\tau}} \frac{\lambda_{\perp} \mathbf{I}}{2\tau} \mathbf{M}^T \right) \quad (101)$$

$$= \mathbf{L} \mathbf{C}^{-1} \mathbf{R} \mathbf{F} \quad (102)$$

643 and

$$\tau \mathbf{L} \left(\frac{d}{dt} \mathbf{C}^{-1} \right) \mathbf{R} = -\tau \mathbf{L} \mathbf{C}^{-1} \left(\frac{d}{dt} \mathbf{C} \right) \mathbf{C}^{-1} \mathbf{R} \quad (103)$$

$$= -\mathbf{L} \mathbf{C}^{-1} \left[\tau \frac{1}{2} \mathbf{Q}(0)^T \mathbf{O} 2 \frac{1}{\tau} e^{2\Lambda \frac{t}{\tau}} \Lambda \Lambda^{-1} \mathbf{O}^T \mathbf{Q}(0) \right. \\ \left. + \tau \frac{1}{2} \mathbf{Q}(0)^T 4 \frac{1}{\tau} \mathbf{M} e^{\lambda \perp \frac{t}{\tau}} \lambda_{\perp} (\lambda_{\perp})^{-1} \mathbf{M}^T \mathbf{Q}(0) \right] \mathbf{C}^{-1} \mathbf{R} \quad (104)$$

$$= -\mathbf{L} \mathbf{C}^{-1} \left[\mathbf{Q}(0)^T \mathbf{O} e^{2\Lambda \frac{t}{\tau}} \mathbf{O}^T \mathbf{Q}(0) + 2 \mathbf{Q}(0)^T \mathbf{M} e^{\lambda \perp \frac{t}{\tau}} \mathbf{M}^T \mathbf{Q}(0) \right] \mathbf{C}^{-1} \mathbf{R} \quad (105)$$

$$= -\mathbf{L} \mathbf{C}^{-1} \left[\mathbf{Q}(0)^T \mathbf{O} e^{\Lambda \frac{t}{\tau}} \mathbf{O}^T \mathbf{O} e^{\Lambda \frac{t}{\tau}} \mathbf{O}^T \mathbf{Q}(0) \right. \\ \left. + 2 \mathbf{Q}(0)^T \mathbf{O} e^{\Lambda \frac{t}{\tau}} \underbrace{\mathbf{O}^T \mathbf{M}}_{\mathbf{0}} e^{\lambda \perp \frac{t}{\tau}} \mathbf{M}^T \mathbf{Q}(0) \right. \\ \left. + 2 \mathbf{Q}(0)^T \mathbf{M} e^{\lambda \perp \frac{t}{\tau}} \underbrace{\mathbf{M}^T \mathbf{O}}_{\mathbf{0}} e^{\Lambda \frac{t}{\tau}} \mathbf{O}^T \mathbf{Q}(0) \right. \\ \left. + 4 \mathbf{Q}(0)^T \mathbf{M} e^{\lambda \perp \frac{t}{\tau}} \mathbf{M}^T \mathbf{M} e^{\lambda \perp \frac{t}{\tau}} \mathbf{M}^T \mathbf{Q}(0) \right] \mathbf{C}^{-1} \mathbf{R} \quad (106)$$

$$= -\mathbf{L} \mathbf{C}^{-1} \mathbf{R} \mathbf{L} \mathbf{C}^{-1} \mathbf{R}. \quad (107)$$

644 Finally, substituting equations 96, 100 and 103 into the left hand side of equation 79 proves equality.
645 \square

646 D Rich-Lazy

647 D.1 Dynamics of the Singular Values

648 **Theorem D.1.** *Under the assumptions of Theorem 2.3 and with a task-aligned initialization given*
649 *by $\mathbf{W}_1(0) = \mathbf{R} \mathbf{S}_1 \tilde{\mathbf{V}}^T$ and $\mathbf{W}_2(0) = \tilde{\mathbf{U}} \mathbf{S}_2 \mathbf{R}^T$, where $\mathbf{R} \in \mathbb{R}^{N_h \times N_h}$ is an orthonormal matrix, then*
650 *the network function is given by the expression $\mathbf{W}_2 \mathbf{W}_1(t) = \tilde{\mathbf{U}} \mathbf{S}(t) \tilde{\mathbf{V}}^T$ where $\mathbf{S}(t) \in \mathbb{R}^{N_h \times N_h}$ is*
651 *a diagonal matrix of singular values with elements $s_{\alpha}(t)$ that evolve according to the equation,*

$$s_{\alpha}(t) = s_{\alpha}(0) + \gamma_{\alpha}(t; \lambda) (\tilde{s}_{\alpha} - s_{\alpha}(0)), \quad (108)$$

652 where \tilde{s}_{α} is the α singular value of $\tilde{\mathbf{S}}$ and $\gamma_{\alpha}(t; \lambda)$ is a λ -dependent monotonic transition function
653 for each singular value that increases from $\gamma_{\alpha}(0; \lambda) = 0$ to $\lim_{t \rightarrow \infty} \gamma_{\alpha}(t; \lambda) = 1$ defined as

$$\gamma_{\alpha}(t; \lambda) = \frac{\tilde{s}_{\lambda, \alpha} s_{\lambda, \alpha} \sinh(2\tilde{s}_{\lambda, \alpha} \frac{t}{\tau}) + \left(\tilde{s}_{\alpha} s_{\alpha} + \frac{\lambda^2}{4} \right) \cosh(2\tilde{s}_{\lambda, \alpha} \frac{t}{\tau}) - \left(\tilde{s}_{\alpha} s_{\alpha} + \frac{\lambda^2}{4} \right)}{\tilde{s}_{\lambda, \alpha} s_{\lambda, \alpha} \sinh(2\tilde{s}_{\lambda, \alpha} \frac{t}{\tau}) + \left(\tilde{s}_{\alpha} s_{\alpha} + \frac{\lambda^2}{4} \right) \cosh(2\tilde{s}_{\lambda, \alpha} \frac{t}{\tau}) + \tilde{s}_{\alpha} (\tilde{s}_{\alpha} - s_{\alpha})}, \quad (109)$$

654 where $\tilde{s}_{\lambda, \alpha} = \sqrt{\tilde{s}_{\alpha}^2 + \frac{\lambda^2}{4}}$, $s_{\lambda, \alpha} = \sqrt{s_{\alpha}(0)^2 + \frac{\lambda^2}{4}}$, and $s_{\alpha} = s_{\alpha}(0)$. We find that under different
655 limits of λ , the transition function converges pointwise to the sigmoidal ($\lambda \rightarrow 0$) and exponential
656 ($\lambda \rightarrow \pm\infty$) transition functions,

$$\gamma_{\alpha}(t; \lambda) \rightarrow \begin{cases} \frac{e^{2\tilde{s}_{\alpha} \frac{t}{\tau}} - 1}{e^{2\tilde{s}_{\alpha} \frac{t}{\tau}} - 1 + \frac{\tilde{s}_{\alpha}}{s_{\alpha}(0)}} & \text{as } \lambda \rightarrow 0, \\ 1 - e^{-|\lambda| \frac{t}{\tau}} & \text{as } \lambda \rightarrow \pm\infty \end{cases}. \quad (110)$$

657 *Proof.* According to Theorem 2.3, the network function is given by the equation

$$\mathbf{W}_2 \mathbf{W}_1(t) = \mathbf{Z}_2(t) \mathbf{A}^{-1}(t) \mathbf{Z}_1^T(t), \quad (111)$$

658 which depends on the variables of the initialization \mathbf{B} and \mathbf{C} . Plugging the expressions for a task-
659 aligned initialization $\mathbf{W}_1(0)$ and $\mathbf{W}_2(0)$ into these variables we get the following simplified expres-
660 sions,

$$\mathbf{B} = \mathbf{R} \underbrace{\left(\mathbf{S}_2(\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}}) + \mathbf{S}_1(\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}}) \right)}_{\mathbf{D}_B}, \quad (112)$$

$$\mathbf{C} = \mathbf{R} \underbrace{\left(\mathbf{S}_2(\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}}) - \mathbf{S}_1(\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}}) \right)}_{\mathbf{D}_C}, \quad (113)$$

661 where we define the diagonal matrices \mathbf{D}_B and \mathbf{D}_C for ease of notation. Using these expressions,
662 we now get the following time-dependent expressions for $\mathbf{Z}_2(t)$, $\mathbf{A}^{-1}(t)$, and $\mathbf{Z}_1(t)$,

$$\mathbf{Z}_1(t) = \frac{1}{2} \tilde{\mathbf{V}} \left((\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}}) e^{\tilde{\mathbf{S}}_\lambda \frac{t}{\tau}} \mathbf{D}_B - (\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}}) e^{-\tilde{\mathbf{S}}_\lambda \frac{t}{\tau}} \mathbf{D}_C \right) \mathbf{R}^T \quad (114)$$

$$\mathbf{Z}_2(t) = \frac{1}{2} \tilde{\mathbf{U}} \left((\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}}) e^{\tilde{\mathbf{S}}_\lambda \frac{t}{\tau}} \mathbf{D}_B + (\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}}) e^{-\tilde{\mathbf{S}}_\lambda \frac{t}{\tau}} \mathbf{D}_C \right) \mathbf{R}^T \quad (115)$$

$$\mathbf{A}(t) = \mathbf{R} \left(\mathbf{I} + \left(\frac{e^{2\tilde{\mathbf{S}}_\lambda \frac{t}{\tau}} - \mathbf{I}}{4\tilde{\mathbf{S}}_\lambda} \right) \mathbf{D}_B^2 - \left(\frac{e^{-2\tilde{\mathbf{S}}_\lambda \frac{t}{\tau}} - \mathbf{I}}{4\tilde{\mathbf{S}}_\lambda} \right) \mathbf{D}_C^2 \right) \mathbf{R}^T \quad (116)$$

663 Plugging these expressions into the expression for the network function, notice that the \mathbf{R} terms
664 cancel each other resulting in following equation

$$\mathbf{W}_2 \mathbf{W}_1(t) = \tilde{\mathbf{U}} \underbrace{\left(\frac{\left((\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}}) e^{\tilde{\mathbf{S}}_\lambda \frac{t}{\tau}} \mathbf{D}_B - (\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}}) e^{-\tilde{\mathbf{S}}_\lambda \frac{t}{\tau}} \mathbf{D}_C \right) \left((\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}}) e^{\tilde{\mathbf{S}}_\lambda \frac{t}{\tau}} \mathbf{D}_B + (\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}}) e^{-\tilde{\mathbf{S}}_\lambda \frac{t}{\tau}} \mathbf{D}_C \right)}{4\mathbf{I} + \left(\frac{e^{2\tilde{\mathbf{S}}_\lambda \frac{t}{\tau}} - \mathbf{I}}{\tilde{\mathbf{S}}_\lambda} \right) \mathbf{D}_B^2 - \left(\frac{e^{-2\tilde{\mathbf{S}}_\lambda \frac{t}{\tau}} - \mathbf{I}}{\tilde{\mathbf{S}}_\lambda} \right) \mathbf{D}_C^2} \right)}_{\mathbf{S}(t)} \tilde{\mathbf{V}}^T, \quad (117)$$

665 Notice that the middle term is simply a product of diagonal matrices. We can factor the numerator
666 of this expressions as,

$$(\tilde{\mathbf{G}}^2 - \tilde{\mathbf{H}}^2 \tilde{\mathbf{G}}^2) e^{2\tilde{\mathbf{S}}_\lambda \frac{t}{\tau}} \mathbf{D}_B^2 + \left((\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}})^2 - (\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}})^2 \right) \mathbf{D}_B \mathbf{D}_C - (\tilde{\mathbf{G}}^2 - \tilde{\mathbf{H}}^2 \tilde{\mathbf{G}}^2) e^{-2\tilde{\mathbf{S}}_\lambda \frac{t}{\tau}} \mathbf{D}_C^2 \quad (118)$$

667 We can further factor this expression as,

$$\tilde{\mathbf{G}}^2 (\mathbf{I} - \tilde{\mathbf{H}}^2) \left(e^{2\tilde{\mathbf{S}}_\lambda \frac{t}{\tau}} \mathbf{D}_B^2 - e^{-2\tilde{\mathbf{S}}_\lambda \frac{t}{\tau}} \mathbf{D}_C^2 \right) - 4\tilde{\mathbf{G}}^2 \tilde{\mathbf{H}} \mathbf{D}_B \mathbf{D}_C. \quad (119)$$

668 Putting it all together we find that $\mathbf{S}(t)$ can be expressed as,

$$\mathbf{S}(t) = \frac{\tilde{\mathbf{G}}^2 (\mathbf{I} - \tilde{\mathbf{H}}^2) \left(e^{2\tilde{\mathbf{S}}_\lambda \frac{t}{\tau}} \mathbf{D}_B^2 - e^{-2\tilde{\mathbf{S}}_\lambda \frac{t}{\tau}} \mathbf{D}_C^2 \right) - 4\tilde{\mathbf{G}}^2 \tilde{\mathbf{H}} \mathbf{D}_B \mathbf{D}_C}{4\mathbf{I} + \left(\frac{e^{2\tilde{\mathbf{S}}_\lambda \frac{t}{\tau}} - \mathbf{I}}{\tilde{\mathbf{S}}_\lambda} \right) \mathbf{D}_B^2 - \left(\frac{e^{-2\tilde{\mathbf{S}}_\lambda \frac{t}{\tau}} - \mathbf{I}}{\tilde{\mathbf{S}}_\lambda} \right) \mathbf{D}_C^2}. \quad (120)$$

669 Now using the relationship between $\tilde{\mathbf{H}}$ and $\tilde{\mathbf{G}}$ we use the following two identities:

$$\tilde{\mathbf{G}}^2 (\mathbf{I} - \tilde{\mathbf{H}}^2) = \frac{\tilde{\mathbf{S}}}{\tilde{\mathbf{S}}_\lambda}, \quad 4\tilde{\mathbf{G}}^2 \tilde{\mathbf{H}} = \frac{\lambda}{\tilde{\mathbf{S}}_\lambda} \quad (121)$$

670 Plugging these identities into the previous expression and multiplying the numerator and denomina-
671 tor by $\tilde{\mathbf{S}}_\lambda$ gives,

$$\mathbf{S}(t) = \frac{\tilde{\mathbf{S}} \left(e^{2\tilde{\mathbf{S}}_\lambda \frac{t}{\tau}} \mathbf{D}_B^2 - e^{-2\tilde{\mathbf{S}}_\lambda \frac{t}{\tau}} \mathbf{D}_C^2 \right) - \lambda \mathbf{D}_B \mathbf{D}_C}{4\tilde{\mathbf{S}}_\lambda + e^{2\tilde{\mathbf{S}}_\lambda \frac{t}{\tau}} \mathbf{D}_B^2 - e^{-2\tilde{\mathbf{S}}_\lambda \frac{t}{\tau}} \mathbf{D}_C^2 + \mathbf{D}_C^2 - \mathbf{D}_B^2}. \quad (122)$$

672 Add and subtract $\tilde{\mathbf{S}} \left(4\tilde{\mathbf{S}}_\lambda + \mathbf{D}_C^2 - \mathbf{D}_B^2 \right)$ from the numerator such that

$$\mathbf{S}(t) = \tilde{\mathbf{S}} - \frac{\tilde{\mathbf{S}} \left(4\tilde{\mathbf{S}}_\lambda + \mathbf{D}_C^2 - \mathbf{D}_B^2 \right) + \lambda \mathbf{D}_B \mathbf{D}_C}{4\tilde{\mathbf{S}}_\lambda + e^{2\tilde{\mathbf{S}}_\lambda \frac{t}{\tau}} \mathbf{D}_B^2 - e^{-2\tilde{\mathbf{S}}_\lambda \frac{t}{\tau}} \mathbf{D}_C^2 + \mathbf{D}_C^2 - \mathbf{D}_B^2}. \quad (123)$$

673 Using the form of D_B and D_C notice the following two identities:

$$D_B D_C = \frac{\lambda}{\tilde{S}_\lambda} \left(\tilde{S} - S_2 S_1 \right), \quad D_C^2 - D_B^2 = -\frac{4}{\tilde{S}_\lambda} \left(\tilde{S} S_2 S_1 + \frac{\lambda^2}{4} \mathbf{I} \right) \quad (124)$$

674 From the second identity we can derive a third identity,

$$4\tilde{S}_\lambda + D_C^2 - D_B^2 = 4 \frac{\tilde{S}}{\tilde{S}_\lambda} \left(\tilde{S} - S_2 S_1 \right) \quad (125)$$

675 Plugging the first and third identities into the numerator for the previous expression gives,

$$S(t) = \tilde{S} - \frac{\frac{(4\tilde{S}^2 + \lambda^2 \mathbf{I})}{\tilde{S}_\lambda} \left(\tilde{S} - S_2 S_1 \right)}{4\tilde{S}_\lambda + e^{2\tilde{S}_\lambda \frac{t}{\tau}} D_B^2 - e^{-2\tilde{S}_\lambda \frac{t}{\tau}} D_C^2 + D_C^2 - D_B^2}. \quad (126)$$

676 Multiply numerator and denominator by $\frac{\tilde{S}_\lambda}{4}$ and simplify terms gives the expression,

$$S(t) = \tilde{S} - \frac{\tilde{S}_\lambda^2}{\tilde{S}_\lambda^2 + \frac{\tilde{S}_\lambda}{4} \left(e^{2\tilde{S}_\lambda \frac{t}{\tau}} D_B^2 - e^{-2\tilde{S}_\lambda \frac{t}{\tau}} D_C^2 \right) - \frac{\tilde{S}_\lambda}{4} (D_B^2 - D_C^2)}{\tilde{S}_\lambda^2 + \frac{\tilde{S}_\lambda}{4} \left(e^{2\tilde{S}_\lambda \frac{t}{\tau}} D_B^2 - e^{-2\tilde{S}_\lambda \frac{t}{\tau}} D_C^2 \right) - \frac{\tilde{S}_\lambda}{4} (D_B^2 - D_C^2)} \left(\tilde{S} - S_2 S_1 \right). \quad (127)$$

677 Thus we have found the transition function,

$$\gamma(t; \lambda) = \frac{\frac{\tilde{S}_\lambda}{4} \left(e^{2\tilde{S}_\lambda \frac{t}{\tau}} D_B^2 - e^{-2\tilde{S}_\lambda \frac{t}{\tau}} D_C^2 \right) + \frac{\tilde{S}_\lambda}{4} (D_C^2 - D_B^2)}{\frac{\tilde{S}_\lambda}{4} \left(e^{2\tilde{S}_\lambda \frac{t}{\tau}} D_B^2 - e^{-2\tilde{S}_\lambda \frac{t}{\tau}} D_C^2 \right) + \frac{\tilde{S}_\lambda}{4} \left(4\tilde{S}_\lambda + D_C^2 - D_B^2 \right)}. \quad (128)$$

678 We will use our previous identities and the definitions of D_B^2 and D_C^2 to simplify this expression.
679 Notice the following identity,

$$\frac{\tilde{S}_\lambda}{4} \left(e^{2\tilde{S}_\lambda \frac{t}{\tau}} D_B^2 - e^{-2\tilde{S}_\lambda \frac{t}{\tau}} D_C^2 \right) = \tilde{S}_\lambda S_\lambda \sinh \left(2\tilde{S}_\lambda \frac{t}{\tau} \right) + \left(\tilde{S} S(0) + \frac{\lambda^2}{4} \mathbf{I} \right) \cosh \left(2\tilde{S}_\lambda \frac{t}{\tau} \right) \quad (129)$$

680 Putting it all together we get

$$\gamma(t; \lambda) = \frac{\tilde{S}_\lambda S_\lambda \sinh \left(2\tilde{S}_\lambda \frac{t}{\tau} \right) + \left(\tilde{S} S(0) + \frac{\lambda^2}{4} \mathbf{I} \right) \cosh \left(2\tilde{S}_\lambda \frac{t}{\tau} \right) - \left(\tilde{S} S(0) + \frac{\lambda^2}{4} \mathbf{I} \right)}{\tilde{S}_\lambda S_\lambda \sinh \left(2\tilde{S}_\lambda \frac{t}{\tau} \right) + \left(\tilde{S} S(0) + \frac{\lambda^2}{4} \mathbf{I} \right) \cosh \left(2\tilde{S}_\lambda \frac{t}{\tau} \right) + \tilde{S} \left(\tilde{S} - S(0) \right)} \quad (130)$$

681 We will now show why under certain limits of λ this expression simplifies to the sigmoidal and
682 exponential dynamics discussed in the previous section.

683 **Sigmoidal dynamics.** When $\lambda = 0$, then $\tilde{S}_\lambda = \tilde{S}$ and $S_\lambda = S(0)$. Notice, that the coefficients for
684 the hyperbolic functions all simplify to $\tilde{S} S(0)$. Using the hyperbolic identity $\sinh(x) + \cosh(x) =$
685 e^x , we can simplify the expression for the transition function to

$$\gamma(t; \lambda) = \frac{\tilde{S} S(0) e^{2\tilde{S} \frac{t}{\tau}} - \tilde{S} S(0)}{\tilde{S} S(0) e^{2\tilde{S} \frac{t}{\tau}} - \tilde{S} S(0) + \tilde{S}^2}. \quad (131)$$

686 Dividing the numerator and denominator by $\tilde{S} S(0)$ gives the final expression.

687 **Exponential dynamics.** In the limit as $\lambda \rightarrow \pm\infty$ the expressions $\tilde{S}_\lambda \rightarrow \frac{|\lambda|}{2}$ and $S_\lambda \rightarrow \frac{|\lambda|}{2}$.
688 Additionally, in these limits because $\frac{\lambda^2}{4} \mathbf{I} \gg \tilde{S} S(0)$ then $\left(\tilde{S} S(0) + \frac{\lambda^2}{4} \mathbf{I} \right) \rightarrow \frac{\lambda^2}{4} \mathbf{I}$. As a result of
689 these simplifications the coefficients for the hyperbolic functions all simplify to $\frac{\lambda^2}{4} \mathbf{I}$. As a result we
690 can again use the hyperbolic identity $\sinh(x) + \cosh(x) = e^x$ to simplify the expression as

$$\gamma(t; \lambda) = \frac{\frac{\lambda^2}{4} e^{|\lambda| \frac{t}{\tau}} - \frac{\lambda^2}{4} \mathbf{I}}{\frac{\lambda^2}{4} e^{|\lambda| \frac{t}{\tau}} + \tilde{S} \left(\tilde{S} - S(0) \right)}. \quad (132)$$

691 Dividing the numerator and denominator by $\frac{\lambda^2}{4}$ results in all terms without a coefficient proportional
692 to λ^2 vanishing, which simplifying further gives the final expression. \square

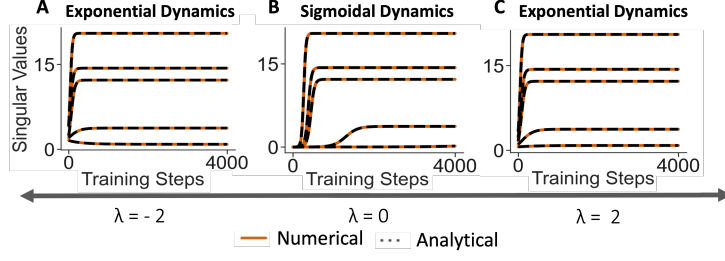


Figure 4: Simulated and analytical dynamics of the singular values of the network function with *relative scale* lambda **A** $\lambda = -2$ **B** $\lambda = 0$ **C** $\lambda = 2$ initialized as described in F.7.

693 D.2 Dynamics of the representation from the Lazy to the Rich Regime

694 The *lazy* and *rich* regimes are defined by the dynamics of the NTK of the network. *Lazy* learning
 695 occurs when the NTK is constant, *rich* learning occurs when it is not. (Farrell et al. (2023b))

696 The NTK intuitively measures the movement of the network representations through training. As
 697 shown in (Braun et al. (2022)), in specific experimental setup, we can calculate the NTK of the
 698 network in terms of the internal representations in a straightforward way:

$$\text{NTK} = \mathbf{I}_{N_o} \otimes \mathbf{X}^T \mathbf{W}_1^T \mathbf{W}_1(t) \mathbf{X} + \mathbf{W}_2 \mathbf{W}_2^T(t) \otimes \mathbf{X}^T \mathbf{X} \quad (133)$$

699 In order to better understand the effect of λ on NTK dynamics, we first prove some theorems in-
 700 volving the Singular Values of the λ -balanced weights, and the representations of a λ -balanced
 701 network.

702 D.2.1 Lambda-balanced singular value

703 **Theorem D.2.** Under a λ -Balanced initialization 2, if the network function $\mathbf{W}_2 \mathbf{W}_1(t) =$
 704 $\mathbf{U}(t) \mathbf{S}(t) \mathbf{V}^T(t)$ is full-rank 4 and we define $\mathcal{S}_\lambda(t) = \sqrt{\mathbf{S}^2(t) + \frac{\lambda^2}{4} \mathbf{I}}$, then we can recover the pa-
 705 rameters $\mathbf{W}_2(t) = \mathbf{U}(t) \mathbf{S}_2(t) \mathbf{R}^T(t)$, $\mathbf{W}_1(t) = \mathbf{R}(t) \mathbf{S}_1(t) \mathbf{V}^T(t)$ up to time-dependent orthogonal
 706 transformation $\mathbf{R}(t)$ of size $N_h \times N_h$, where

$$\mathbf{S}_1(t) = \left(\left(\mathcal{S}_\lambda(t) - \frac{\lambda \mathbf{I}}{2} \right)^{\frac{1}{2}} \quad 0_{\max(0, N_i - N_o)} \right) \quad \mathbf{S}_2(t) = \left(\left(\mathcal{S}_\lambda(t) + \frac{\lambda \mathbf{I}}{2} \right)^{\frac{1}{2}} \quad 0_{\max(0, N_o - N_i)} \right) \quad (134)$$

707 *Proof.* We prove the case $N_i \leq N_o$ and $N_h = \min(N_i, N_o)$. The proof for $N_o \leq N_i$ follows the
 708 same structure. Let $\mathbf{U} \mathbf{S} \mathbf{V}^T = \mathbf{W}_2(t) \mathbf{W}_1(t)$ be the Singular Value Decomposition of the product
 709 of the weights at training step t . We will use $\mathbf{W}_2 = \mathbf{W}_2(t)$, $\mathbf{W}_1 = \mathbf{W}_1(t)$ as a shorthand.

710 By properties of Singular Value Decomposition, we can write $\mathbf{W}_2 = \mathbf{U} \mathbf{S}_2 \mathbf{R}^T$, $\mathbf{W}_1 = \mathbf{R} \mathbf{S}_1 \mathbf{V}^T$,
 711 where \mathbf{R} is an orthonormal matrix and $\mathbf{S}_2, \mathbf{S}_1$ are diagonal (possibly rectangular) matrices.

712 The Balanced property states that $\mathbf{W}_2^T \mathbf{W}_2 - \mathbf{W}_1 \mathbf{W}_1^T = \lambda \mathbf{I}$. We know this holds for any t
 713 since this is a conserved quantity in linear networks.

714

715 Hence

$$\mathbf{R} \mathbf{S}_2^T \mathbf{S}_2 \mathbf{R}^T - \mathbf{R} \mathbf{S}_1 \mathbf{S}_1 \mathbf{R}^T = \lambda \mathbf{I} \quad (135)$$

$$\mathbf{S}_2^T \mathbf{S}_2 - \mathbf{S}_1 \mathbf{S}_1 = \lambda \mathbf{I} \quad (136)$$

716 The matrices $\mathbf{S}_1, \mathbf{S}_2$, have shapes $(N_h, N_i), (N_o, N_h)$ respectively. We introduce the diagonal ma-
 717 trices $\hat{\mathbf{S}}_1$ of shape (N_h, N_i) , $\hat{\mathbf{S}}_2$ of shape (N_i, N_h) such that the zero matrix has size $(N_o - N_i, N_h)$
 718 :
 719 :
 720 :

$$\mathbf{S}_1 = (\hat{\mathbf{S}}_1), \quad \mathbf{S}_2 = \begin{pmatrix} \hat{\mathbf{S}}_2 \\ 0 \end{pmatrix} \quad (137)$$

721 Hence

$$\mathbf{S}_2^T \mathbf{S}_2 - \mathbf{S}_1 \mathbf{S}_1 = \lambda \mathbf{I} \quad (138)$$

722 From the equation above and the fact that $\hat{\mathbf{S}}_1 \hat{\mathbf{S}}_2 = \mathbf{S}$ we derive that:

$$\hat{\mathbf{S}}_2 = \left(\frac{\sqrt{\lambda^2 \mathbf{I} + 4\mathbf{S}^2} + \lambda \mathbf{I}}{2} \right)^{\frac{1}{2}}, \quad \hat{\mathbf{S}}_1 = \left(\frac{\sqrt{\lambda^2 \mathbf{I} + 4\mathbf{S}^2} - \lambda \mathbf{I}}{2} \right)^{\frac{1}{2}}, \quad (139)$$

723 Hence

$$\mathbf{W}_2 = \mathbf{U} \left(\begin{array}{c} \left(\frac{\sqrt{\lambda^2 \mathbf{I} + 4\mathbf{S}^2} + \lambda \mathbf{I}}{2} \right)^{\frac{1}{2}} \\ 0_{\max(0, N_o - N_i)} \end{array} \right), \mathbf{R}^T, \quad \mathbf{W}_1 = \mathbf{R} \left(\begin{array}{c} \left(\frac{\sqrt{\lambda^2 \mathbf{I} + 4\mathbf{S}^2} - \lambda \mathbf{I}}{2} \right)^{\frac{1}{2}} \\ 0_{\max(0, N_i - N_o)} \end{array} \right) \mathbf{V}^T \quad (140)$$

724 \square

725 D.2.2 Convergence proof

726 With our solution, $\mathbf{Q}\mathbf{Q}^T(t)$, which captures the temporal dynamics of the similarity between hidden
 727 layer activations, we can analyze the network's internal representations in relation to the task. This
 728 allows us to determine whether the network adopts a *rich* or *lazy* representation, depending on the
 729 value of λ . Consider a λ -Balanced network training on data $\Sigma^{yx} = \tilde{\mathbf{U}}\tilde{\mathbf{S}}\tilde{\mathbf{V}}^T$. We assume that the
 730 convergence is toward global minima and \mathbf{B} is invertible

731 **Theorem D.3.** *Under the assumptions of Theorem C.5, the network function converges to $\tilde{\mathbf{U}}\tilde{\mathbf{S}}\tilde{\mathbf{V}}^T$
 732 and acquires the internal representation, that is $\mathbf{W}_1^T \mathbf{W}_1 = \tilde{\mathbf{V}}\tilde{\mathbf{S}}_1^2 \tilde{\mathbf{V}}^T$ and $\mathbf{W}_2 \mathbf{W}_2^T = \tilde{\mathbf{U}}\tilde{\mathbf{S}}_2^2 \tilde{\mathbf{U}}^T$*

733 *Proof.* As training time increases, all terms including a matrix exponential with negative exponent
 734 in Equation 70 vanish to zero, as $\mathbf{S}_\lambda = \tilde{\mathbf{S}}_\lambda$ is a diagonal matrix with entries larger zero

735 As training time increases, all terms in the equations vanish to zero. Terms in Equation 70 decay as

$$\lim_{t \rightarrow \infty} e^{-\sqrt{\tilde{\mathbf{S}}^2 + \frac{\lambda^2 \mathbf{I}}{4}} \frac{t}{\tau}} = \mathbf{0}. \quad (141)$$

736 and

$$\lim_{t \rightarrow \infty} e^{\lambda \frac{t}{\tau}} e^{-\sqrt{\tilde{\mathbf{S}}^2 + \frac{\lambda^2 \mathbf{I}}{4}} \frac{t}{\tau}} = \mathbf{0}. \quad (142)$$

737 where $\tilde{\mathbf{S}}_\lambda = \tilde{\mathbf{S}}_\lambda$ is a diagonal matrix with entries larger zero

738 Therefore, in the temporal limit, eq. 70 reduces to

$$\lim_{t \rightarrow \infty} \mathbf{Q}\mathbf{Q}^T(t) = \lim_{t \rightarrow \infty} \begin{bmatrix} \mathbf{W}_1^T \mathbf{W}_1(t) & \mathbf{W}_1^T \mathbf{W}_2^T(t) \\ \mathbf{W}_2 \mathbf{W}_1(t) & \mathbf{W}_2 \mathbf{W}_2(t) \end{bmatrix} \quad (143)$$

$$= \begin{bmatrix} \tilde{\mathbf{V}}(\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}}) \\ \tilde{\mathbf{U}}(\tilde{\mathbf{H}}\tilde{\mathbf{G}} + \tilde{\mathbf{G}}) \end{bmatrix} \left[\tilde{\mathbf{S}}_\lambda^{-1} \right]^{-1} \left[(\tilde{\mathbf{V}}(\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}}))^T \quad (\tilde{\mathbf{U}}(\tilde{\mathbf{H}}\tilde{\mathbf{G}} + \tilde{\mathbf{G}}))^T \right] \quad (144)$$

$$= \begin{bmatrix} \tilde{\mathbf{V}}(\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}})\tilde{\mathbf{S}}_\lambda(\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}})^T \tilde{\mathbf{V}}^T & \tilde{\mathbf{V}}(\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}})\tilde{\mathbf{S}}_\lambda(\tilde{\mathbf{H}}\tilde{\mathbf{G}} + \tilde{\mathbf{G}})^T \tilde{\mathbf{U}}^T \\ \tilde{\mathbf{U}}(\tilde{\mathbf{H}}\tilde{\mathbf{G}} + \tilde{\mathbf{G}})\tilde{\mathbf{S}}_\lambda(\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}})^T \tilde{\mathbf{V}}^T & \tilde{\mathbf{U}}(\tilde{\mathbf{H}}\tilde{\mathbf{G}} + \tilde{\mathbf{G}})\tilde{\mathbf{S}}_\lambda(\tilde{\mathbf{H}}\tilde{\mathbf{G}} + \tilde{\mathbf{G}})^T \tilde{\mathbf{U}}^T \end{bmatrix}. \quad (145)$$

$$(\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}})\tilde{\mathbf{S}}_\lambda(\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}}) = \frac{\mathbf{S}_\lambda(1 - \tilde{\mathbf{H}}^2)}{1 + \tilde{\mathbf{H}}^2} = \tilde{\mathbf{S}} \quad (146)$$

$$\tilde{\mathbf{S}}_\lambda(\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}})^2 = \frac{\tilde{\mathbf{S}}_\lambda(1 + \tilde{\mathbf{H}}^2)}{1 + \tilde{\mathbf{H}}^2} - \frac{\tilde{\mathbf{S}}_\lambda(2\tilde{\mathbf{H}})}{1 + \tilde{\mathbf{H}}^2} = \frac{\sqrt{4\tilde{\mathbf{S}}^2 + \lambda^2 \mathbf{I}} - \lambda \mathbf{I}}{2} \quad (147)$$

$$\tilde{S}_\lambda(\tilde{G} + \tilde{H}\tilde{G})^2 = \frac{\tilde{S}_\lambda(1 + \tilde{H}^2)}{1 + \tilde{H}^2} + \frac{\tilde{S}_\lambda(2\tilde{H})}{1 + \tilde{H}^2} = \frac{\sqrt{4\tilde{S}^2 + \lambda^2\mathbf{I}} + \lambda\mathbf{I}}{2} \quad (148)$$

$$\lim_{t \rightarrow \infty} \mathbf{Q}\mathbf{Q}^T(t) = \lim_{t \rightarrow \infty} \begin{bmatrix} \mathbf{W}_1^T \mathbf{W}_1(t) & \mathbf{W}_1^T \mathbf{W}_2^T(t) \\ \mathbf{W}_2^T \mathbf{W}_1(t) & \mathbf{W}_2^T \mathbf{W}_2(t) \end{bmatrix} \quad (149)$$

$$= \begin{bmatrix} \tilde{\mathbf{V}}\tilde{S}_1^2\tilde{\mathbf{V}}^T & \tilde{\mathbf{V}}\tilde{S}\tilde{\mathbf{U}}^T \\ \tilde{\mathbf{U}}\tilde{S}\tilde{\mathbf{V}}^T & \tilde{\mathbf{U}}\tilde{S}_2^2\tilde{\mathbf{U}}^T \end{bmatrix}. \quad (150)$$

739

□

740 D.2.3 Representation in the limit

741 **Theorem D.4.** Under the assumptions of Theorem C.5, training on data $\Sigma^{yx} = \tilde{\mathbf{U}}\tilde{\mathbf{S}}\tilde{\mathbf{V}}^T$, as $\lambda \rightarrow \infty$
742 the representation tends to

$$\mathbf{W}_2\mathbf{W}_2^T = \tilde{\mathbf{U}} \begin{pmatrix} \lambda\mathbf{I} & 0_{\max(0, N_o - N_i)} \\ 0_{\max(0, N_o - N_i)} & 0 \end{pmatrix} \tilde{\mathbf{U}}^T \quad \mathbf{W}_1^T\mathbf{W}_1 = \frac{1}{\lambda} \tilde{\mathbf{V}} \begin{pmatrix} \tilde{S}^2 & 0_{\max(0, N_i - N_o)} \\ 0_{\max(0, N_i - N_o)} & 0 \end{pmatrix} \tilde{\mathbf{V}}^T$$

743 As $\lambda \rightarrow -\infty$

$$\mathbf{W}_2\mathbf{W}_2^T = -\frac{1}{\lambda} \tilde{\mathbf{U}} \begin{pmatrix} \tilde{S}^2 & 0_{\max(0, N_o - N_i)} \\ 0_{\max(0, N_o - N_i)} & 0 \end{pmatrix} \tilde{\mathbf{U}}^T, \quad \mathbf{W}_1^T\mathbf{W}_1 = \tilde{\mathbf{V}} \begin{pmatrix} -\lambda\mathbf{I} & 0_{\max(0, N_i - N_o)} \\ 0_{\max(0, N_i - N_o)} & 0 \end{pmatrix} \tilde{\mathbf{V}}^T$$

744 As $\lambda \rightarrow -\infty$

$$\mathbf{W}_2\mathbf{W}_2^T = -\frac{1}{\lambda} \tilde{\mathbf{U}} \begin{pmatrix} \tilde{S}^2 & 0_{\max(0, N_o - N_i)} \\ 0_{\max(0, N_o - N_i)} & 0 \end{pmatrix} \tilde{\mathbf{U}}^T, \quad \mathbf{W}_1^T\mathbf{W}_1 = \tilde{\mathbf{V}} \begin{pmatrix} -\lambda\mathbf{I} & 0_{\max(0, N_i - N_o)} \\ 0_{\max(0, N_i - N_o)} & 0 \end{pmatrix} \tilde{\mathbf{V}}^T$$

745 *Proof.* We start from the representation derived in D.3 and using the Taylor expansion of $f(x) =$
746 $\sqrt{1 + x^2}$, we compute

$$\frac{\sqrt{\lambda^2\mathbf{I} + 4\tilde{S}^2} + \lambda\mathbf{I}}{2} = \frac{|\lambda|\sqrt{1 + \left(\frac{2\tilde{S}}{\lambda}\right)^2} + \lambda\mathbf{I}}{2} \quad (151)$$

$$\frac{|\lambda|\left(1 + \left(\frac{2\tilde{S}}{\lambda}\right)^2 + O(\lambda^{-4})\right) + \lambda\mathbf{I}}{2} = \frac{|\lambda| + \lambda}{2} + \frac{\tilde{S}^2}{|\lambda|} + O(\lambda^{-3}) \quad (152)$$

747 Hence

$$\lim_{\lambda \rightarrow \infty} \frac{\sqrt{\lambda^2\mathbf{I} + 4\tilde{S}^2} + \lambda\mathbf{I}}{2} = \lambda\mathbf{I}, \quad \lim_{\lambda \rightarrow -\infty} \frac{\sqrt{\lambda^2\mathbf{I} + 4\tilde{S}^2} + \lambda\mathbf{I}}{2} = \frac{\tilde{S}^2}{|\lambda|} = -\frac{\tilde{S}^2}{\lambda} \quad (153)$$

748 Similarly,

$$\frac{\sqrt{\lambda^2\mathbf{I} + 4\tilde{S}^2} - \lambda\mathbf{I}}{2} = \frac{|\lambda| - \lambda}{2} + \frac{\tilde{S}^2}{|\lambda|} + O(\lambda^{-3}) \quad (154)$$

$$\lim_{\lambda \rightarrow \infty} \frac{\sqrt{\lambda^2\mathbf{I} + 4\tilde{S}^2} - \lambda\mathbf{I}}{2} = \frac{\tilde{S}^2}{\lambda}, \quad \lim_{\lambda \rightarrow -\infty} \frac{\sqrt{\lambda^2\mathbf{I} + 4\tilde{S}^2} - \lambda\mathbf{I}}{2} = \frac{\tilde{S}^2}{|\lambda|} = -\lambda\mathbf{I} \quad (155)$$

749 Since $\tilde{\mathbf{U}}, \tilde{\mathbf{V}}$ are independent of λ :

$$\lim_{\lambda \rightarrow \pm\infty} \mathbf{W}_2 \mathbf{W}_2^T = \tilde{\mathbf{U}} \left(\lim_{\lambda \rightarrow \pm\infty} \mathbf{S}_2 \right) \tilde{\mathbf{U}}^T \quad (156)$$

$$\lim_{\lambda \rightarrow \pm\infty} \mathbf{W}_1^T \mathbf{W}_1 = \tilde{\mathbf{V}} \left(\lim_{\lambda \rightarrow \pm\infty} \mathbf{S}_1 \right) \tilde{\mathbf{V}}^T \quad (157)$$

750

□

751 As $|\lambda| \rightarrow \infty$, one of the network representations approaches a scaled identity matrix, while the
 752 other tends toward zero. Intuitively, this suggests that the representations shift less and less as $|\lambda|$
 753 increases. Next, we demonstrate that the NTK becomes progressively less variable as $|\lambda|$ grows and
 754 ultimately converges to zero.

755 D.2.4 NTK movement

756 Relationship between λ and the NTK of the network

757 **Theorem D.5.** *Under the assumptions of Theorem C.5, consider a linear network training on data*
 758 $\Sigma^{yx} = \tilde{\mathbf{U}} \tilde{\mathbf{S}} \tilde{\mathbf{V}}^T$. *At any arbitrary training time $t \geq 0$, let $\mathbf{W}_2(t) \mathbf{W}_1(t) = \mathbf{U}^* \mathbf{S}^* \mathbf{V}^{*T}$. Then,*

759 1. For any $\lambda \in \mathbf{R}$:

$$\begin{aligned} \text{NTK}(0) &= \mathbf{I}_{N_o} \otimes \mathbf{X}^T \mathbf{V} \begin{pmatrix} \frac{\sqrt{\lambda^2 \mathbf{I} + 4\mathbf{S}^{*2}} - \lambda \mathbf{I}}{2} & 0 \\ 0 & 0 \end{pmatrix} \mathbf{V}^T \mathbf{X} \\ &\quad + \mathbf{U} \begin{pmatrix} \frac{\sqrt{\lambda^2 \mathbf{I} + 4\mathbf{S}^{*2}} + \lambda \mathbf{I}}{2} & 0 \\ 0 & 0 \end{pmatrix} \mathbf{U}^T \otimes \mathbf{X}^T \mathbf{X} \end{aligned} \quad (158)$$

$$\begin{aligned} \text{NTK}(t) &= \mathbf{I}_{N_o} \otimes \mathbf{X}^T \mathbf{V}^* \begin{pmatrix} \frac{\sqrt{\lambda^2 \mathbf{I} + 4\mathbf{S}^{*2}} - \lambda \mathbf{I}}{2} & 0 \\ 0 & 0 \end{pmatrix} \mathbf{V}^{*T} \\ &\quad + \mathbf{U}^* \begin{pmatrix} \frac{\sqrt{\lambda^2 \mathbf{I} + 4\mathbf{S}^{*2}} + \lambda \mathbf{I}}{2} & 0 \\ 0 & 0 \end{pmatrix} \mathbf{U}^{*T} \otimes \mathbf{X}^T \mathbf{X} \end{aligned} \quad (159)$$

760 2. As $\lambda \rightarrow \infty$:

$$\text{NTK}(t) - \text{NTK}(0) \rightarrow \frac{1}{\lambda} \left(\mathbf{I}_{N_o} \otimes \mathbf{X}^T \mathbf{V}^* \tilde{\mathbf{S}}^{*2} \mathbf{V}^{*T} \mathbf{X} - \mathbf{I}_{N_o} \otimes \mathbf{X}^T \mathbf{V} \tilde{\mathbf{S}}^2 \mathbf{V}^T \mathbf{X} \right) \rightarrow 0 \quad (160)$$

761 3. As $\lambda \rightarrow -\infty$:

$$\text{NTK}(t) - \text{NTK}(0) \rightarrow \frac{1}{\lambda} \left(\mathbf{U} \tilde{\mathbf{S}}^2 \mathbf{U}^T \otimes \mathbf{X}^T \mathbf{X} - \mathbf{U}^* \tilde{\mathbf{S}}^{*2} \mathbf{U}^{*T} \otimes \mathbf{X}^T \mathbf{X} \right) \rightarrow 0 \quad (161)$$

762 **Proof.** Follows by substituting the expressions for the network representations in terms of λ from
 763 (Braun et al. (2022))’s expression for the NTK of a linear network. Similarly, follows from substi-
 764 tuting the limit expressions for the network representations and the fact that the Kronecker product
 765 is linear in both arguments. □

766 The theorem above demonstrates that as $|\lambda| \rightarrow \infty$, the NTK of a λ -Balanced network remains
 767 constant. This indicates that the network operates in the *lazy* regime throughout all training steps.
 768 This finding is significant as it highlights the impact of weight initialization on learning regimes.

769 **D.3 Representation robustness and sensitivity to noise**

770 As derived in (Braun et al., 2024), the expected mean squared error under additive, independent and
771 identically distributed input noise with mean $\mu = 0$ and variance $\sigma_{\mathbf{x}}^2$ is

$$\left\langle \frac{1}{2P} \sum_{i=1}^P \|\mathbf{W}_2 \mathbf{W}_1 (\mathbf{x}_x + \xi_i) - \mathbf{y}_i\|_2^2 \right\rangle_{\xi_{\mathbf{x}}} = \sigma_{\mathbf{x}}^2 \|\mathbf{W}_2 \mathbf{W}_1\|_F^2 + c, \quad (162)$$

772 where $c = \frac{1}{2} \text{Tr}(\tilde{\Sigma}^{yy}) - \frac{1}{2} \text{Tr}(\tilde{\Sigma}^{yx} \tilde{\Sigma}^{yxT})$ is a noise independent constant that only depends on
773 the statistics of the training data. In Theorem D.3 we show that the network function converges to
774 $\tilde{\mathbf{U}} \tilde{\mathbf{S}} \tilde{\mathbf{V}}^T$ and therefore

$$\begin{aligned} \sigma_{\mathbf{x}}^2 \|\mathbf{W}_2 \mathbf{W}_1\|_F^2 &= \sigma_{\mathbf{x}}^2 \|\tilde{\mathbf{U}} \tilde{\mathbf{S}} \tilde{\mathbf{V}}^T\|_F^2 \\ &= \sigma_{\mathbf{x}}^2 \|\tilde{\mathbf{S}}\|_F^2 \\ &= \sigma_{\mathbf{x}}^2 \sum_{i=1}^{N_h} \tilde{\mathbf{S}}_i^2 \end{aligned} \quad (163)$$

775 As derived in (Braun et al., 2024), under the assumption of whitened inputs (Assumption 1), in the
776 case of additive parameter noise with $\mu = 0$ and variance $\sigma_{\mathbf{W}}^2$, the expected mean squared error is

$$\begin{aligned} &\left\langle \frac{1}{2P} \sum_{i=1}^P \|(\mathbf{W}_2 + \xi \mathbf{w}_2) (\mathbf{W}_1 + \xi \mathbf{w}_1) \mathbf{x}_i - \mathbf{y}_i\|_2^2 \right\rangle_{\xi_{\mathbf{w}_1}, \xi_{\mathbf{w}_2}} \\ &= \frac{1}{2} N_i \sigma_{\mathbf{W}}^2 \|\mathbf{W}_2\|_F^2 + \frac{1}{2} N_o \sigma_{\mathbf{W}}^2 \|\mathbf{W}_1\|_F^2 + \frac{1}{2} N_i N_h N_o \sigma^4 + c. \end{aligned} \quad (164)$$

777 Using Theorem D.3, we have

$$\begin{aligned} \|\mathbf{W}_1\|_F^2 &= \text{Tr}(\mathbf{W}_1^T \mathbf{W}_1) \\ &= \text{Tr} \left(\frac{\sqrt{\lambda^2 \mathbf{I} + 4\tilde{\mathbf{S}}^2} + \lambda \mathbf{I}}{2} \right) \\ &= \frac{1}{2} \left(\sum_{i=1}^{N_h} \sqrt{\lambda^2 + 4\tilde{\mathbf{S}}_i^2} + \lambda \right) \end{aligned} \quad (165)$$

778 and

$$\begin{aligned} \|\mathbf{W}_2\|_F^2 &= \text{Tr}(\mathbf{W}_2 \mathbf{W}_2^T) \\ &= \text{Tr} \left(\frac{\sqrt{\lambda^2 \mathbf{I} + 4\tilde{\mathbf{S}}^2} - \lambda \mathbf{I}}{2} \right) \\ &= \frac{1}{2} \left(\sum_{i=1}^{N_h} \sqrt{\lambda^2 + 4\tilde{\mathbf{S}}_i^2} - \lambda \right). \end{aligned} \quad (166)$$

779 To find the λ that minimises the expected loss, we substitute the equations for the norms, take the
780 partial derivative with respect to λ and set it to zero

$$\begin{aligned} &\frac{\partial \langle \mathcal{L} \rangle_{\xi_{\mathbf{w}_1}, \xi_{\mathbf{w}_2}}}{\partial \lambda} \stackrel{!}{=} 0 \\ &\Leftrightarrow \frac{1}{4} N_i \sigma_{\mathbf{W}}^2 \frac{\partial}{\partial \lambda} \left(\sum_{i=1}^{N_h} \sqrt{\lambda^2 + 4\tilde{\mathbf{S}}_i^2} - \lambda \right) + \frac{1}{4} N_o \sigma_{\mathbf{W}}^2 \frac{\partial}{\partial \lambda} \left(\sum_{i=1}^{N_h} \sqrt{\lambda^2 + 4\tilde{\mathbf{S}}_i^2} + \lambda \right) = 0 \\ &\Leftrightarrow N_i \sum_{i=1}^{N_h} \frac{\lambda}{\sqrt{\lambda^2 + 4\tilde{\mathbf{S}}_i^2}} - N_i N_h + N_o \sum_{i=1}^{N_h} \frac{\lambda}{\sqrt{\lambda^2 + 4\tilde{\mathbf{S}}_i^2}} + N_o N_h = 0 \\ &\Leftrightarrow \sum_{i=1}^{N_h} \frac{\lambda}{\sqrt{\lambda^2 + 4\tilde{\mathbf{S}}_i^2}} = N_h \frac{N_i - N_o}{N_i + N_o}. \end{aligned} \quad (167)$$

781 It follows, that under the assumption that $N_i = N_o$, the equation reduces to

$$\sum_{i=1}^{N_h} \frac{\lambda}{\sqrt{\lambda^2 + 4\tilde{\mathbf{S}}_i^2}} = 0. \quad (168)$$

782 We note, that the denominator is always positive and therefore, that the left-hand side of the equation
 783 is always larger zero for any $\lambda > 0$, and smaller than zero for any $\lambda < 0$. The equation is therefore
 784 only solved for $\lambda = 0$.

785 **D.4 Effect of the architecture from the lazy to the Rich Regime**

786 **Theorem D.6.** *Under the conditions of Theorem C.5, when $\lambda_{\perp} > 0$, the network enters a regime*
 787 *referred to as the delayed-rich phase. In this phase, the learning rate is determined by two competing*
 788 *exponential factors:*

$$e^{\lambda_{\perp} \frac{t}{\tau}} e^{-\sqrt{\tilde{\mathbf{S}}^2 + \frac{\lambda^2}{4} \mathbf{I}} \frac{t}{\tau}}$$

789 and

$$e^{-\sqrt{\tilde{\mathbf{S}}^2 + \frac{\lambda^2}{4} \mathbf{I}} \frac{t}{\tau}}.$$

790 As λ increases, different parts of the network exhibit distinct learning behaviors: some components
 791 adapt quickly and converge exponentially with lambda, while others are constrained by the singular
 792 values of the network, resulting in slower adaptation.

793 *Proof.* The solution to Theorem C.5 is governed by two time-dependent terms:

$$e^{-\sqrt{\tilde{\mathbf{S}}^2 + \frac{\lambda^2}{4} \mathbf{I}} \frac{t}{\tau}} \quad \text{and} \quad e^{\lambda_{\perp} \frac{t}{\tau}} e^{-\sqrt{\tilde{\mathbf{S}}^2 + \frac{\lambda^2}{4} \mathbf{I}} \frac{t}{\tau}}.$$

794 The first term exhibits exponential decay with rate λ , approaching zero as time progresses:

$$\lim_{t \rightarrow \infty} e^{-\sqrt{\tilde{\mathbf{S}}^2 + \frac{\lambda^2}{4} \mathbf{I}} \frac{t}{\tau}} = \mathbf{0}.$$

795 The second term also decays, but at a rate governed by the singular values $\tilde{\mathbf{S}}$, as λ tends to infinity:

$$\lim_{t \rightarrow \infty} e^{\lambda_{\perp} \frac{t}{\tau}} e^{-\sqrt{\tilde{\mathbf{S}}^2 + \frac{\lambda^2}{4} \mathbf{I}} \frac{t}{\tau}} = \mathbf{0}.$$

796 Since

$$\lambda_{\perp} - \sqrt{\tilde{\mathbf{S}}^2 + \frac{\lambda^2}{4} \mathbf{I}} > 0,$$

797 we have

$$\lim_{\lambda \rightarrow \infty} \left(\lambda_{\perp} - \sqrt{\tilde{\mathbf{S}}^2 + \frac{\lambda^2}{4} \mathbf{I}} \right) = \tilde{\mathbf{S}}.$$

798 Thus, as λ increases, the convergence rate slows for certain parts of the network (those governed by
 799 larger singular values), while other components continue to learn more quickly. This explains the
 800 delay observed in the delayed-rich regime. \square

801 **E Appendix: Application**

802 **E.1 Appendix: Continual Learning**

We build upon the derivation presented in Braun et al. (2022) to incorporate the dynamics of continual learning throughout the entire learning trajectory. Utilizing the assumption of whitened inputs,

the entire batch loss for the i th task is

$$\begin{aligned}
\mathcal{L}_i(\mathcal{T}_j) &= \frac{1}{2P} \|\mathbf{W}_2 \mathbf{W}_1 \mathbf{X}_i - \mathbf{Y}_i\|_F^2 \\
&= \frac{1}{2P} \text{Tr} \left((\mathbf{W}_2 \mathbf{W}_1 \mathbf{X}_i - \mathbf{Y}_i) (\mathbf{W}_2 \mathbf{W}_1 \mathbf{X}_i - \mathbf{Y}_i)^T \right) \\
&= \frac{1}{2P} \text{Tr} (\mathbf{W}_2 \mathbf{W}_1 \mathbf{X}_i \mathbf{X}_i^T (\mathbf{W}_2 \mathbf{W}_1)^T) - \frac{1}{P} \text{Tr} (\mathbf{W}_2 \mathbf{W}_1 \mathbf{X}_i \mathbf{Y}_i^T) + \frac{1}{2P} \text{Tr} (\mathbf{Y}_i \mathbf{Y}_i^T) \\
&= \frac{1}{2} \text{Tr} (\mathbf{W}_2 \mathbf{W}_1 (\mathbf{W}_2 \mathbf{W}_1)^T) - \text{Tr} (\mathbf{W}_2 \mathbf{W}_1 \tilde{\Sigma}_i^{yx^T}) + \frac{1}{2} \text{Tr} (\tilde{\Sigma}_i^{yy}) \\
&= \frac{1}{2} \text{Tr} \left((\mathbf{W}_2 \mathbf{W}_1 - \tilde{\Sigma}_i^{yx}) (\mathbf{W}_2 \mathbf{W}_1 - \tilde{\Sigma}_i^{yx})^T - \tilde{\Sigma}_i^{yx} \tilde{\Sigma}_i^{yx^T} \right) + \frac{1}{2} (\tilde{\Sigma}_i^{yy}) \\
&= \frac{1}{2} \left\| \mathbf{W}_2 \mathbf{W}_1 - \tilde{\Sigma}_i^{yx} \right\|_F^2 - \underbrace{\frac{1}{2} \text{Tr} (\tilde{\Sigma}_i^{yx} \tilde{\Sigma}_i^{yx^T}) + \frac{1}{2} (\tilde{\Sigma}_i^{yy})}_c.
\end{aligned}$$

803 Hence, the extent of forgetting, denoted as \mathcal{F} for task \mathcal{T}_i during training on task \mathcal{T}_k subsequent to
804 training the network on task \mathcal{T}_j , specifically, the relative change in loss, is entirely dictated by the
805 similarity structure among tasks.

$$\begin{aligned}
\mathcal{F}_i(\mathcal{T}_j, \mathcal{T}_k) &= \mathcal{L}_i(\mathcal{T}_k) - \mathcal{L}_i(\mathcal{T}_j) \\
&= \frac{1}{2} \left\| \tilde{\Sigma}_k^{yx} - \tilde{\Sigma}_i^{yx} \right\|_F^2 + c - \frac{1}{2} \left\| \mathbf{W}_2 \mathbf{W}_1 - \tilde{\Sigma}_i^{yx} \right\|_F^2 - c \\
&= \frac{1}{2} \left(\left\| \tilde{\Sigma}_k^{yx} - \tilde{\Sigma}_i^{yx} \right\|_F^2 - \left\| \mathbf{W}_2 \mathbf{W}_1 - \tilde{\Sigma}_i^{yx} \right\|_F^2 \right).
\end{aligned}$$

806 It is important to note that the amount of forgetting is a function of the weight trajectories. Therefore,
807 we have analytical solutions for trajectories of forgetting as well.

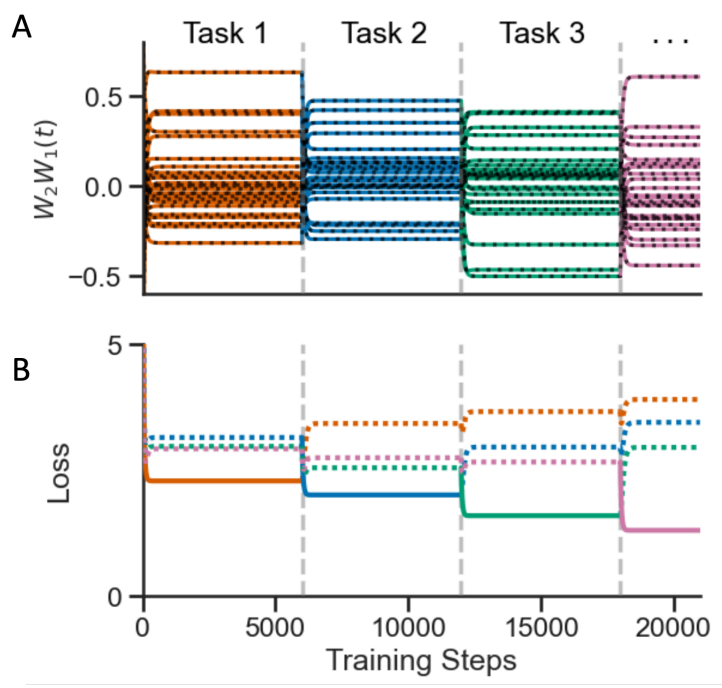


Figure 5: Continual learning. **A** Top: Network training from small zero-balanced weights across a sequence of tasks (colored lines represent simulations, and black dotted lines represent analytical results). Bottom: Evaluation loss for the tasks in the sequence (dotted lines) while training on the current task (solid lines). As the network optimizes its function on the current task, the loss on previously learned tasks increases.

808 Figure. E.1 panel was generated by training a linear network with $N_i = 5$, $N_h = 10$, $N_o = 6$
809 subsequently on four different random regression tasks with $N = 25$. The learning rate was $\eta =$
810 0.05 and the initial weights were small ($\sigma = 0.0001$).

811 **E.2 Appendix: Reversal Learning**

812 As first introduced in Braun et al. (2022), in the following discussion, we assume that the input and
813 output dimensions are equal. We denote the i -th columns of the left and right singular vectors as \mathbf{u}_i ,
814 $\tilde{\mathbf{u}}_i$, and \mathbf{v}_i , $\tilde{\mathbf{v}}_i$, respectively.

815 Reversal learning occurs when both the task and the initial network function share the same left and
816 right singular vectors, i.e., $\mathbf{U} = \tilde{\mathbf{U}}$ and $\mathbf{V} = \tilde{\mathbf{V}}$, with the exception of one or more columns of the
817 left singular vectors, where the direction is reversed: $-\mathbf{u}_i = \tilde{\mathbf{u}}_i$.

818 It is important to note that if a reversal occurs in the right singular vectors, such that $-\mathbf{v}_i = \tilde{\mathbf{v}}_i$, this
819 can be equivalently represented as a reversal in the left singular vectors, as the signs of the right and
820 left singular vectors are interchangeable.

821 In the reversal learning setting, both $\mathbf{B} = \mathbf{S}_2 \tilde{\mathbf{U}}^T \tilde{\mathbf{U}} (\tilde{\mathbf{G}} + \tilde{\mathbf{H}} \tilde{\mathbf{G}}) + \mathbf{S}_1 \mathbf{V}^T \tilde{\mathbf{V}} (\tilde{\mathbf{G}} - \tilde{\mathbf{H}} \tilde{\mathbf{G}})$ and
822 $\mathbf{C} = \mathbf{S}_2 \tilde{\mathbf{U}}^T \tilde{\mathbf{U}} (\tilde{\mathbf{G}} - \tilde{\mathbf{H}} \tilde{\mathbf{G}}) - \mathbf{S}_1 \mathbf{V}^T \tilde{\mathbf{V}} (\tilde{\mathbf{G}} + \tilde{\mathbf{H}} \tilde{\mathbf{G}})$ are diagonal matrices.

823

824 In the case where lambda is zero, the same argument given in Braun et al. (2022) follows, the
825 diagonal entries of \mathbf{C} are zero if the singular vectors are aligned and non zero if they are reversed.
826 Similarly, diagonal entries of \mathbf{B} are non-zero if the singular vectors are aligned and zero if they are
827 reversed. Therefore, in the case of reversal learning, \mathbf{B} is a diagonal matrix with 0 values and thus
828 is not invertible. As a consequence, the learning dynamics cannot be described by Equation 49.
829 However, as \mathbf{B} and \mathbf{C} are diagonal matrices, the learning dynamics simplify. Let \mathbf{b}_i , \mathbf{c}_i , \mathbf{s}_i and $\tilde{\mathbf{s}}_i$

830 denote the i -th diagonal entry of \mathbf{B} , \mathbf{C} , \mathbf{S} and $\tilde{\mathbf{S}}$ respectively, then the network dynamics can be
 831 rewritten as

$$\mathbf{W}_2 \mathbf{W}_1(t) = \frac{1}{2} \tilde{\mathbf{U}} \left[(\tilde{\mathbf{G}} + \tilde{\mathbf{H}} \tilde{\mathbf{G}}) e^{\tilde{\mathbf{S}}_\lambda \frac{t}{\tau}} \mathbf{B}^T + (\tilde{\mathbf{G}} - \tilde{\mathbf{H}} \tilde{\mathbf{G}}) e^{-\tilde{\mathbf{S}}_\lambda \frac{t}{\tau}} \mathbf{C}^T \right] \left[\mathbf{S}_\lambda^{-1} + \frac{1}{4} \mathbf{B} \left(e^{2\tilde{\mathbf{S}}_\lambda \frac{t}{\tau}} - \mathbf{I} \right) \tilde{\mathbf{S}}_\lambda^{-1} \mathbf{B}^T - \frac{1}{4} \mathbf{C} \left(e^{-2\tilde{\mathbf{S}}_\lambda \frac{t}{\tau}} - \mathbf{I} \right) \tilde{\mathbf{S}}_\lambda^{-1} \mathbf{C}^T \right]^{-1} \quad (169)$$

$$\frac{1}{2} \left((\tilde{\mathbf{G}} - \tilde{\mathbf{H}} \tilde{\mathbf{G}}) e^{\tilde{\mathbf{S}}_\lambda \frac{t}{\tau}} \mathbf{B} - (\tilde{\mathbf{G}} + \tilde{\mathbf{H}} \tilde{\mathbf{G}}) e^{-\tilde{\mathbf{S}}_\lambda \frac{t}{\tau}} \mathbf{C} \right) \tilde{\mathbf{V}}^T$$

$$= \sum_{i=1}^{N_i} \frac{\mathbf{b}_i^2 e^{2\tilde{\mathbf{s}}_{\lambda i} \frac{t}{\tau}} - \mathbf{c}_i^2 e^{-2\tilde{\mathbf{s}}_{\lambda i} \frac{t}{\tau}}}{4\mathbf{s}_{\lambda i}^{-1} + \mathbf{b}_i^2 e^{2\tilde{\mathbf{s}}_{\lambda i} \frac{t}{\tau}} \tilde{\mathbf{s}}_{\lambda i}^{-1} - \mathbf{b}_i^2 \tilde{\mathbf{s}}_{\lambda i}^{-1} - \mathbf{c}_i^2 e^{-2\tilde{\mathbf{s}}_{\lambda i} \frac{t}{\tau}} \tilde{\mathbf{s}}_{\lambda i}^{-1} + \mathbf{c}_i^2 \tilde{\mathbf{s}}_{\lambda i}^{-1}} \tilde{\mathbf{u}}_i \tilde{\mathbf{v}}_i^T \quad (170)$$

$$= \sum_{i=1}^{N_i} \frac{\mathbf{s}_{\lambda i} \mathbf{b}_i^2 \tilde{\mathbf{s}}_{\lambda i} - \mathbf{s}_{\lambda i} \mathbf{c}_i^2 \tilde{\mathbf{s}}_{\lambda i} e^{-4\tilde{\mathbf{s}}_{\lambda i} \frac{t}{\tau}}}{4\tilde{\mathbf{s}}_{\lambda i} e^{-2\tilde{\mathbf{s}}_{\lambda i} \frac{t}{\tau}} + \mathbf{s}_{\lambda i} \mathbf{b}_i^2 \left(1 - e^{-2\tilde{\mathbf{s}}_{\lambda i} \frac{t}{\tau}} \right) + \mathbf{s}_{\lambda i} \mathbf{c}_i^2 \left(e^{-2\tilde{\mathbf{s}}_{\lambda i} \frac{t}{\tau}} - e^{-4\tilde{\mathbf{s}}_{\lambda i} \frac{t}{\tau}} \right)} \tilde{\mathbf{u}}_i \tilde{\mathbf{v}}_i^T \quad (171)$$

832 It follows, that in the reversal learning case, i.e. $\mathbf{b} = 0$, for each reversed singular vector, the
 833 dynamics vanish to zero

$$\lim_{t \rightarrow \infty} \frac{-\mathbf{s}_{\lambda i} \mathbf{c}_i^2 \tilde{\mathbf{s}}_{\lambda i} e^{-4\tilde{\mathbf{s}}_{\lambda i} \frac{t}{\tau}}}{4\tilde{\mathbf{s}}_{\lambda i} e^{-2\tilde{\mathbf{s}}_{\lambda i} \frac{t}{\tau}} + \mathbf{s}_i \mathbf{c}_i^2 \left(e^{-2\tilde{\mathbf{s}}_{\lambda i} \frac{t}{\tau}} - e^{-4\tilde{\mathbf{s}}_{\lambda i} \frac{t}{\tau}} \right)} \tilde{\mathbf{u}}_i \tilde{\mathbf{v}}_i^T = 0. \quad (172)$$

834 Analytically, the learning dynamics are initialized on and remain along the separatrix of a saddle
 835 point until the corresponding singular value of the network function decreases to zero and stays
 836 there, indicating convergence to the saddle point. In numerical simulations, however, the learning
 837 dynamics can escape the saddle points due to the imprecision of floating-point arithmetic. Despite
 838 this, numerical optimization still experiences significant delays, as escaping the saddle point is time-
 839 consuming Lee et al. (2022). In contrast, when the singular vectors are aligned ($\mathbf{c} = 0$), the equation
 840 governing temporal dynamics, as described in Saxe et al. (2014), is recovered. Under these con-
 841 ditions, training succeeds, with the singular value of the network function converging to its target
 842 value.

$$\lim_{t \rightarrow \infty} \sum_{i=1}^{N_i} \frac{\mathbf{s}_{\lambda i} \mathbf{b}_i^2 \tilde{\mathbf{s}}_{\lambda i}}{4\tilde{\mathbf{s}}_{\lambda i} e^{-2\tilde{\mathbf{s}}_{\lambda i} \frac{t}{\tau}} + \mathbf{s}_{\lambda i} \mathbf{b}_i^2 \left(1 - e^{-2\tilde{\mathbf{s}}_{\lambda i} \frac{t}{\tau}} \right)} \tilde{\mathbf{u}}_i \tilde{\mathbf{v}}_i^T = \frac{\mathbf{s}_{\lambda i} \mathbf{b}_i^2 \tilde{\mathbf{s}}_{\lambda i}}{\mathbf{s}_{\lambda i} \mathbf{b}_i^2} \tilde{\mathbf{u}}_i \tilde{\mathbf{v}}_i^T \quad (173)$$

$$= \tilde{\mathbf{s}}_{\lambda i} \tilde{\mathbf{u}}_i \tilde{\mathbf{v}}_i^T. \quad (174)$$

843 In summary, in the case of aligned singular vectors, the learning dynamics can be described by
 844 the convergence of singular values. However in the case of reversal learning, analytically, training
 845 does not succeed. In simulations, the learning dynamics escape the saddle point due to numerical
 846 imprecision, but the learning dynamics are catastrophically slowed in the vicinity of the saddle point
 847 as shown in figure E.2 .

848 In the case where λ is non-zero, the diagonal of \mathbf{C} are also non-zero; this is true regardless of
 849 whether they are reversed or aligned. Similarly, the diagonal entries of \mathbf{B} remain non-zero whether
 850 the singular vectors are aligned or reversed. Therefore, in the case of reversal learning, \mathbf{B} is a
 851 diagonal matrix with elements that are zero. In figure E.2

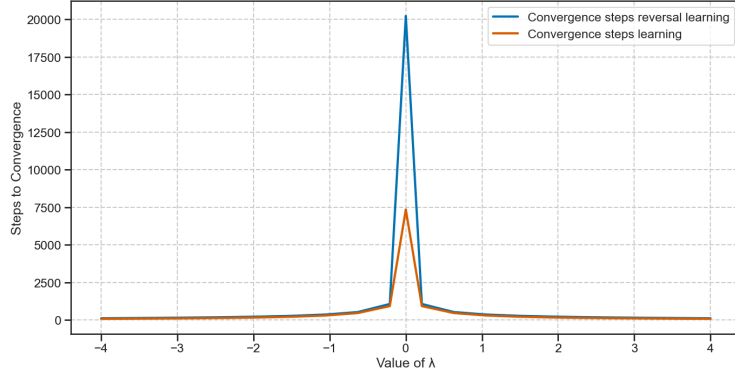


Figure 6: Plot showing the steps to convergence for two tasks: (1) the reversal learning task and (2) a randomly sampled continual learning task across a range of λ values. The reversal learning task exhibits catastrophic slowing at $\lambda = 0$.

852 E.3 Appendix: Generalization and structured learning

853 We study how the representations learned for different λ initializations impact generalization of
 854 properties of the data. To do this, we consider the case where a new feature is associated to a
 855 learned item in a dataset and how this new feature may then be related to other items based on prior
 856 knowledge. In particular, we first train each network (for different values of $-10 \leq \lambda \leq 10$) on
 857 the hierarchical semantic learning task in Section 3 and then add a new feature (e.g., ‘eats worms’)
 858 to a single item (e.g., the goldfish) (Fig. E.3A), correspondingly increasing the output dimension
 859 to represent the novel feature. In order to learn the new feature without affecting prior knowledge,
 860 we append a randomly initialized row to \mathbf{W}_2 and train it on the single item with the new feature,
 861 while keeping the rest of the network frozen. Thus, we only change the weights from the hidden
 862 layer to the new feature which may produce different behavior depending on how the hidden layer
 863 representations vary based on λ . After training on the new feature-item association, we query the
 864 network with the rest of the data to observe how the new feature is associated with the other items.
 865 We find that as λ increases positively, the network better transfers the hierarchy such that it projects
 866 the feature onto items based on their distance to the trained item (Fig. E.3B,C). For example, after
 867 learning that a goldfish eats worms, the network can extrapolate the hierarchy to infer that another
 868 fish, or birds, may also eat worms; instead, plants are not likely to eat worms. Alternatively, as λ
 869 becomes more negative, the network ceases to infer any hierarchical structure and only learns to map
 870 the new feature to the single item trained on. In this case, after learning that a goldfish eats worms,
 871 the network does not infer that other fish, birds, or plants may also eat worms.

872 Interestingly, this setting highlights how asymmetries in the representations yielded by different λ
 873 can actually benefit transfer and generalization. This can be shown by observing that the learning
 874 of a new feature association only depends on the first layer \mathbf{W}_1 . Let $\hat{\mathbf{y}}_f$ denote the vector of the
 875 representation of the new feature f across items i in the dataset. Additionally, let $\mathbf{w}_2^{(f)T}$ be the new
 876 row of weights appended to \mathbf{W}_2 which map the hidden layer to the new feature. Following Saxe
 877 et al. (2019b), if $\mathbf{w}_2^{(f)T}$ is initialized with small random weights and trained on item $\tilde{\mathbf{H}}_i$, it will
 878 converge to

$$\mathbf{w}_2^{(f)T} = \tilde{\mathbf{H}}_i^T \mathbf{W}_1^T / \|\mathbf{W}_1 \tilde{\mathbf{H}}_i\|_2^2 \quad (175)$$

$$\hat{\mathbf{y}}_f = (\tilde{\mathbf{H}}_i^T \mathbf{W}_1^T \mathbf{W}_1 \tilde{\mathbf{H}}) / \|\mathbf{W}_1 \tilde{\mathbf{H}}_i\|_2^2 \quad (176)$$

879 From this we can see that differences in the representations of the new feature across items $\hat{\mathbf{y}}_f$ across
 880 λ are only influenced by \mathbf{W}_1 .

881 In the case of the rich learning regime where $\lambda = 0$, the semantic relationship between features
 882 and items is distributed across both layers. Instead, when $\lambda > 0$, the second layer \mathbf{W}_2 exhibits
 883 *lazy* learning, yielding an output representation $\mathbf{W}_2 \mathbf{W}_2^T$ of a weighted identity matrix. However,
 884 the first layer \mathbf{W}_1 still learns a *rich* representation of the hierarchy, albeit at a smaller scaling.
 885 Furthermore, rather than distributing this learning across both layers, in the $\lambda > 0$ case, all learning

886 of the hierarchy occurs in the first layer, allowing it to more readily transfer this structure to the
 887 learning of a new feature (which only depends on the first layer). Thus, in this case, the ‘shallowing’
 888 of the network into the first layer is actually beneficial. Finally, we can also observe the opposite
 889 case when $\lambda < 0$. Here, *rich* learning happens in the second layer, while the first layer is *lazy* and
 890 learns to represent a weighted identity matrix. As such, these networks do not learn to transfer the
 891 hierarchy of different items to the new feature.

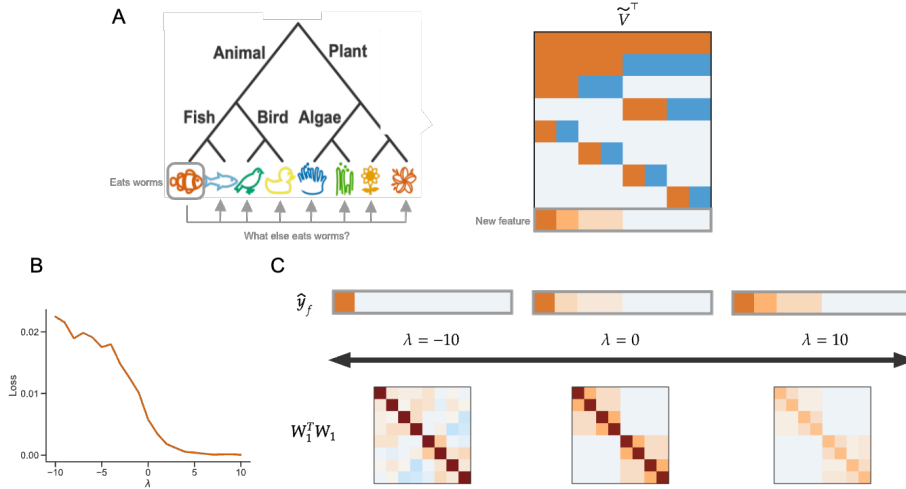


Figure 7: Transfer learning for different λ . **A** A new feature (such as ‘eats worms’) is introduced to the dataset after training on the hierarchical semantic learning task (Section 3). A randomly initialized row is added to \mathbf{W}_2 and trained on a single item with the new feature (for example, the goldfish), with the rest of the network frozen. The network is then tested on the transfer of the new feature to other items, such that items closer to the goldfish in the hierarchy are more likely to have the same feature. **B** The generalization loss on the untrained items with the new feature decreases as λ increases. **C** As λ increases positively, networks better transfer the hierarchical structure of the data to the representation of the new feature.

892 F Implementation and Simulations

893 The details of the simulation studies are described as follows. Specifically, N_i , N_h , and N_o represent
 894 the dimensions of the input, hidden layer, and output (target), respectively. The total number of
 895 training samples is denoted by N , and the learning rate is defined as $\eta = \frac{1}{\tau}$.

896 F.1 Lambda-balanced weight initialization

897 In practice, to initialize the network with lambda-balanced weights, we use Algorithm F.1. In this
 898 algorithm, α serves as a scaling factor that controls the variance of the weights, allowing for adjust-
 899 ments between smaller and larger weight initializations.

900 F.2 Tasks

901 In the following, we describe the different tasks that are used throughout the simulation studies.

902 F.2.1 Random regression task

903 In the random regression task, the inputs $\mathbf{X} \in \mathbb{R}^{N_i \times N}$ are generated from a standard normal dis-
 904 tribution, $\mathbf{X} \sim \mathcal{N}(\mu = 0, \sigma = 1)$. The input data \mathbf{X} is then whitened to satisfy $\frac{1}{N} \mathbf{X} \mathbf{X}^T = \mathbf{I}$.
 905 The target values $\mathbf{Y} \in \mathbb{R}^{N_o \times N}$ are independently sampled from a normal distribution with variance
 906 scaled according to the number of output nodes, $\mathbf{Y} \sim \mathcal{N}(\mu = 0, \alpha = \frac{1}{N_o})$. Consequently, the

Algorithm 1 Get λ -balanced

```
1: function GET_LAMBDA_BALANCED( $\lambda$ ,  $in\_dim$ ,  $hidden\_dim$ ,  $out\_dim$ ,  $\sigma = 1$ )
2:   if  $out\_dim > in\_dim$  and  $\lambda < 0$  then
3:     raise Exception('Lambda must be positive if out_dim  $\geq$  in_dim')
4:   end if
5:   if  $in\_dim > out\_dim$  and  $\lambda > 0$  then
6:     raise Exception('Lambda must be positive if in_dim  $\geq$  out_dim')
7:   end if
8:   if  $hidden\_dim < \min(in\_dim, out\_dim)$  then
9:     raise Exception('Network cannot be bottlenecked')
10:  end if
11:  if  $hidden\_dim > \max(in\_dim, out\_dim)$  and  $\lambda \neq 0$  then
12:    raise Exception('hidden_dim cannot be the largest dimension if lambda is not 0')
13:  end if
14:   $W_1 \leftarrow \sigma \cdot \text{random normal matrix}(hidden\_dim, in\_dim)$ 
15:   $W_2 \leftarrow \sigma \cdot \text{random normal matrix}(out\_dim, hidden\_dim)$ 
16:   $[U, S, Vt] \leftarrow \text{SVD}(W_2 \cdot W_1)$ 
17:   $R \leftarrow \text{random orthonormal matrix}(hidden\_dim)$ 
18:   $S2_{equal\_dim} \leftarrow \sqrt{(\sqrt{\lambda^2 + 4 \cdot S^2} + \lambda) / 2}$ 
19:   $S1_{equal\_dim} \leftarrow \sqrt{(\sqrt{\lambda^2 + 4 \cdot S^2} - \lambda) / 2}$ 
20:  if  $out\_dim > in\_dim$  then
21:     $S2 \leftarrow \begin{bmatrix} S2_{equal\_dim} & & 0 \\ & & \\ 0 & & 0_{hidden\_dim - in\_dim} \end{bmatrix}$ 
22:     $S1 \leftarrow \begin{bmatrix} S1_{equal\_dim} \\ & \\ 0 & & \end{bmatrix}$ 
23:  else if  $in\_dim > out\_dim$  then
24:     $S1 \leftarrow \begin{bmatrix} S1_{equal\_dim} & & 0 \\ & & \\ 0 & & 0_{hidden\_dim - out\_dim} \end{bmatrix}$ 
25:     $S2 \leftarrow \begin{bmatrix} S2_{equal\_dim} & & \\ & & \\ & & 0 \end{bmatrix}$ 
26:  end if
27:   $init\_W_2 \leftarrow U \cdot S2 \cdot R^T$ 
28:   $init\_W_1 \leftarrow R \cdot S1 \cdot Vt$ 
29:  return ( $init\_W_1, init\_W_2$ )
30: end function
```

907 network inputs and target values are uncorrelated Gaussian noise, implying that a linear solution
908 may not always exist.

909 F.2.2 Semantic hierarchy

910 We use the same task as in Braun et al. (2022) and modify it to match the theoretical dynamics.
911 The modification ensures that the inputs are whitened. In the semantic hierarchy task, input items
912 are represented as one-hot vectors, i.e., $\mathbf{X} = \frac{1}{8}$. The corresponding target vectors, \mathbf{y}_i , encode the
913 item's position within the hierarchical tree. Specifically, a value of 1 indicates that the item is a left
914 child of a node, -1 denotes a right child, and 0 indicates that the item is not a child of that node.
915 For example, consider the blue fish: it is a blue fish, a left child of the root node, a left child of the
916 animal node, not part of the plant branch, a right child of the fish node, and not part of the bird,
917 algae, or flower branches, resulting in the label $[1, 1, 1, 0, -1, 0, 0, 0]$. The labels for all objects in
918 the semantic tree, as shown in Figure 2 A, are given by:

$$\mathbf{Y} = 8 * \begin{bmatrix} 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 & -1 & -1 & -1 & -1 \\ 1 & 1 & -1 & -1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 1 & -1 & -1 \\ 1 & -1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & -1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & -1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & -1 \end{bmatrix}. \quad (177)$$

919 The singular value decomposition (SVD) of the corresponding correlation matrix, $\tilde{\Sigma}^{yx}$, is not unique
 920 due to identical singular values: the first two, the third and fourth, and the last four values are the
 921 same. To align the numerical and analytical solutions, this permutation invariance is addressed by
 922 adding a small perturbation to each column \mathbf{y}_i , for $i \in 1, \dots, N$, of the labels:

$$\mathbf{y}_i = \mathbf{y}_i \cdot \left(1 + \frac{0.1}{i}\right), \quad (178)$$

923 resulting in singular values that are nearly, but not exactly, identical.

924 F.3 Figure 1

925 Panels B illustrates three simulations conducted on the same task with varying initial λ -balanced
 926 weights respectively $\lambda = -2, \lambda = 0, \lambda = 2$. The regression task parameters were set with ($\sigma =$
 927 $\sqrt{10}$). The network architecture consisted of $N_i = 3, N_h = 2, N_o = 2$, with a learning rate of
 928 $\eta = 0.0002$. The batch size is $N = 10$. The zero-balanced weights are initialized with variance
 929 $\sigma = 0.00001$. The lambda-balanced network are initialized with $sigma_{xy} = \sqrt{1}$ of a random
 930 regression task with same architecture.

931 On Panel C, we plot the ballancedness $\mathbf{W}_2(0)^T \mathbf{W}_2(0) - \mathbf{W}_1(0) \mathbf{W}_1(0)^T$ for a two layer network
 932 initialised with Lecun initialization with dimension $N_i = 40, N_h = 120, N_o = 250$

933 F.4 Figure 2

934 Panel A, B, C illustrates three simulations conducted on the same task with varying initial λ -balanced
 935 weights respectively $\lambda = -2, \lambda = 0, \lambda = 2$ according to the initialization scheme described in F.7.
 936 The regression task parameters were set with ($\sigma = \sqrt{10}$). The network architecture consisted of
 937 $N_i = 3, N_h = 2, N_o = 2$ with a learning rate of $\eta = 0.0002$. The batch size is $N = 10$. The
 938 zero-balanced weights are initialized with variance $\sigma = 0.00001$. The lambda-balanced network are
 939 initialized with $sigma_{xy} = \sqrt{1}$ of a random regression task with same architecture.

940 F.5 Figure 3

941 Panel A, B, C illustrates three simulations conducted on the same task with varying initial λ -balanced
 942 weights respectively $\lambda = -2, \lambda = 0, \lambda = 2$ according to the initialization scheme described in F.7.
 943 The regression task parameters were set with ($\sigma = \sqrt{12}$). The network architecture consisted of
 944 $N_i = 3, N_h = 3, N_o = 3$ with a learning rate of $\eta = 0.0002$. The batch size is $N = 5$. The
 945 zero-balanced weights are initialized with variance $\sigma = 0.0009$. The lambda-balanced network are
 946 initialized with $sigma_{xy} = \sqrt{12}$ of a random regression task with same architecture.

947 F.6 Figure 4

948 In Panel A presents a semantic learning task with the SVD of the input-output correlation matrix
 949 of the task. U and V represent the singular vectors, and S contains the singular values. This
 950 decomposition allows us to compute the respective RSMs as USU^T for the input and VSV^T for
 951 the output task. The rows and columns in the SVD and RSMs are ordered identically to the items in
 952 the hierarchical tree.

953 The results in Panel B display simulation outcomes, while Panel C presents theoretical input and
 954 output representation matrices at convergence for a network trained on the semantic task described
 955 in Braun et al. (2022); Saxe et al. (2013). These matrices are generated using varying initial λ -
 956 balanced weights set at $\lambda = -2$, $\lambda = 0$, and $\lambda = 2$, following the initialization scheme outlined
 957 in F.7. The network architecture includes $N_i = 8$, $N_h = 8$, and $N_o = 8$ with a learning rate
 958 of $\eta = 0.001$ and a batch size of $N = 8$. Zero-balanced weights are initialized with a variance
 959 of $\sigma = 0.00001$, while λ -balanced networks are initialized with $\sigma_{xy} = \sqrt{1}$ based on a random
 960 regression task with the same architecture.

961 Panel D illustrates results from running the same task and network configuration but initialized with
 962 randomly large weights having a variance of $\sigma = 1$.

963 In panel E, we trained a two-layer linear network with $N_i = N_h = N_o = 4$ on a random regression
 964 task for $\lambda \in [-5, -4, -3, -2, -1, 0, 1, 2, 3, 4, 5]$ to convergence. Subsequently, we added Gaussian
 965 noise with $\mu = 0, \sigma \in [0, 0.5, 1]$ to the inputs (top panel) or synaptic weights (bottom panel) and
 966 calculated the expected mean squared error.

967 F.7 Figure 5

968 Panel A illustrates schematic representations of the network architectures considered: from left to
 969 right, a funnel network ($N_i = 4, N_h = 2, N_o = 2$), a square network ($N_i = 4, N_h = 4, N_o = 4$),
 970 and an inverted-funnel network ($N_i = 2, N_h = 2, N_o = 4$).

971 Panel B shows the Neural Tangent Kernel (NTK) distance from initialization, as defined in Fort et al.
 972 (2020), across the three architectures shown schematically. The kernel distance is calculated as:

$$S(t) = 1 - \frac{\langle K_0, K_t \rangle}{\|K_0\|_F \|K_t\|_F}.$$

973 The simulations conducted on the same task with eleven varying initial λ -balanced weights in
 974 $[-9, 9]$. The regression task parameters were set with ($\sigma = \sqrt{3}$). The task has batch size $N = 10$.
 975 The network has with a learning rate of $\eta = 0.01$. The lambda-balanced network are initialized with
 976 $\sigma_{xy} = \sqrt{1}$ of a random regression task.

977 Panel C shows the Neural Tangent Kernel (NTK) distance from initialization for the funnel archi-
 978 tectures shown schematically with dimensions $N_i = 3, N_h = 2$, and $N_o = 2$. The simulations
 979 conducted on the same task with twenty one varying initial λ -balanced weights in $[-9, 9]$. The re-
 980 gression task parameters were set with ($\sigma = \sqrt{3}$). The task has batch size $N = 30$. The network
 981 has with a learning rate of $\eta = 0.002$. The lambda-balanced network are initialized with $\sigma_{xy} = \sqrt{1}$
 982 of a random regression task.