

PAPER: ML 2023

Exact learning dynamics of deep linear networks with prior knowledge*

Clémentine C J Dominé^{1,6}, Lukas Braun^{2,6},
James E Fitzgerald³ and Andrew M Saxe^{1,4,5,**}

¹ Gatsby Computational Neuroscience Unit, University College London, London, United Kingdom

² Department of Experimental Psychology, University of Oxford, Oxford, United Kingdom

³ Howard Hughes Medical Institute, Janelia Research Campus, Ashburn, Virginia, VA, United States of America

⁴ Sainsbury Wellcome Centre, University College London, London, United Kingdom

⁵ CIFAR Azrieli Global Scholar, CIFAR, Toronto, Canada

E-mail: a.saxe@ucl.ac.uk

Received 9 June 2023

Accepted for publication 13 September 2023

Published 15 November 2023



Online at stacks.iop.org/JSTAT/2023/114004
<https://doi.org/10.1088/1742-5468/ad01b8>

Abstract. Learning in deep neural networks is known to depend critically on the knowledge embedded in the initial network weights. However, few theoretical results have precisely linked prior knowledge to learning dynamics. Here we derive exact solutions to the dynamics of learning with rich prior knowledge in deep linear networks by generalising Fukumizu’s matrix Riccati solution (Fukumizu 1998 *Gen* 1 1E–03). We obtain explicit expressions for the evolving network function, hidden representational similarity, and neural tangent kernel over training for a broad class of initialisations and tasks. The expressions reveal a class

*This article is an updated version of: Braun L, Dominé C, Fitzgerald J and Saxe A 2022 Exact learning dynamics of deep linear networks with prior knowledge *Advances in Neural Information Processing Systems* vol 35, ed Koyejo S, Mohamed S, Agarwal A, Belgrave D, Cho K and Oh A (Curran Associates, Inc.) pp 6615–29

⁶First authors, random order.

** Author to whom any correspondence should be addressed.



Original Content from this work may be used under the terms of the [Creative Commons Attribution 4.0 licence](https://creativecommons.org/licenses/by/4.0/). Any further distribution of this work must maintain attribution to the author(s) and the title of the work, journal citation and DOI.

of task-independent initialisations that radically alter learning dynamics from slow non-linear dynamics to fast exponential trajectories while converging to a global optimum with identical representational similarity, dissociating learning trajectories from the structure of initial internal representations. We characterise how network weights dynamically align with task structure, rigorously justifying why previous solutions successfully described learning from small initial weights without incorporating their fine-scale structure. Finally, we discuss the implications of these findings for continual learning, reversal learning and learning of structured knowledge. Taken together, our results provide a mathematical toolkit for understanding the impact of prior knowledge on deep learning.

Keywords: deep learning, learning theory, machine learning

Contents

1. Introduction	3
1.1. Contributions	4
1.2. Related work	5
2. Preliminaries and setting	6
3. Exact learning dynamics with prior knowledge	8
4. Rich and lazy learning regimes and generalisation	9
5. Decoupling dynamics	11
6. Applications	13
7. Discussion	16
Acknowledgment	17
Appendix A. Fukumizu approach	17
Appendix B. Network's internal representations	20
B.1. Representational similarity analysis	20
B.2. Finite-width neural tangent kernel	20
Appendix C. Exact learning dynamics with prior knowledge	22
C.1. Proof of theorem 3.1	22
C.1.1. Unequal input-output dimension.	22
C.1.2. Equal input-output dimension.	24
C.2. Derivation of the exact learning dynamics	25
C.2.1. Inverse and matrix exponential of \mathbf{F}	26
C.3. Proof of theorem 3.2: Limiting behaviour	31
C.4. Dynamics of $\mathbf{Q}(\mathbf{t})$	31

Appendix D. Rich and lazy learning regimes and generalisation	32
Appendix E. Decoupling dynamics	33
E.1. Proof for theorem 5.1	33
E.2. Solution for 2×2 dynamics	35
E.3. Off-Diagonal decoupling dynamics	37
E.4. On-diagonal dynamics and the effect of initialisation variance	38
Appendix F. Continual learning	39
Appendix G. Revising structured knowledge	40
G.1. Reversal learning dynamics	40
G.2. Exact learning dynamics in shallow networks	42
Appendix H. Simulations	43
H.1. Zero-balanced weight initialisation	43
H.2. Tasks	44
H.2.1. Random regression task	44
H.2.2. Teacher-student task	44
H.2.3. Semantic hierarchy	44
H.2.4. Colour hierarchy	45
H.3. Figure 1	45
H.4. Figure 2	45
H.5. Figure 3	45
H.6. Figure 4	45
H.7. Figure 5	46
H.8. Figure 6	46
References	46

1. Introduction

A hallmark of human learning is our exquisite sensitivity to prior knowledge: what we already know affects how we subsequently learn (Carey 1985). For instance, having learned about the attributes of nine animals, we may learn about the tenth more quickly (McClelland *et al* 1995, Murphy 2004, McClelland 2013, Flesch *et al* 2018). In machine learning, the impact of prior knowledge on learning is evident in a range of paradigms including reversal learning (Erdeniz and Atalay 2010), transfer learning (Taylor and Stone 2009, Thrun and Pratt 2012, Lampinen and Ganguli 2018, Gerace *et al* 2022), continual learning (Kirkpatrick *et al* 2017, Zenke *et al* 2017, Parisi *et al* 2019), curriculum learning (Bengio *et al* 2009), and meta learning (Javed and White 2019). One form of prior knowledge in deep networks is the initial network state, which is known to strongly impact learning dynamics (Saxe *et al* 2014, Pennington *et al* 2017,

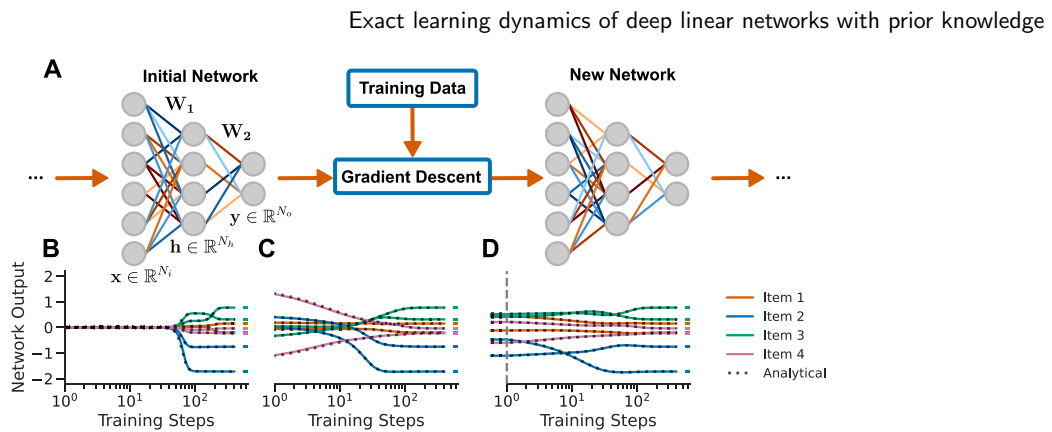


Figure 1. Learning with prior knowledge. (A) In our setting, a deep linear network with N_i input, N_h hidden and N_o output neurons is trained from a particular initialisation using gradient descent. (B)–(D) Network output for an example task over training time when starting from (B) small random weights, (C) large random weights, and (D) the weights of a previously learned task. The dynamics depend in detail on the initialisation. Solid lines indicate simulations, dotted lines indicate the analytical solutions we derive in this work.

Bahri *et al* 2020). Even random initial weights of different variance can yield qualitative shifts in network behaviour between the *lazy* and *rich* regimes (Chizat *et al* 2019), imparting distinct inductive biases on the learning process. More broadly, rich representations such as those obtained through pretraining provide empirically fertile inductive biases for subsequent fine-tuning (Raghu *et al* 2019). Yet while the importance of prior knowledge to learning is clear, our theoretical understanding remains limited, and fundamental questions remain about the implicit inductive biases of neural networks trained from structured initial weights. A better understanding of the impact of initialisation on gradient-based learning may lead to improved pretraining schemes and illuminate pathologies like catastrophic forgetting in continual learning (McCloskey and Cohen 1989).

Here, we address this gap by deriving exact solutions to the dynamics of learning in deep linear networks as a function of network initialisation, revealing an intricate and systematic dependence. We consider the setting depicted in figure 1(A), where a network is trained with standard gradient descent from a potentially complex initialisation. When trained on the same task, different initialisations can radically change the network’s learning trajectory (figures 1(B)–(D)). Our approach, based on a matrix Riccati formalism (Fukumizu 1998), provides explicit analytical expressions for the network output over time (figures 1(B)–(D) dotted). While simple, deep linear networks have a non-convex loss landscape and have been shown to recapitulate several features of nonlinear deep networks while retaining mathematical tractability.

1.1. Contributions

- We derive an explicit solution for the gradient flow of the network function, internal representational similarity, and finite-width neural tangent kernel (NTK) of over- and

under-complete two-layer deep linear networks for a rich class of initial conditions (section 3).

- We characterise a set of random initial network states that exhibit fast, exponential learning dynamics and yet converge to *rich* neural representations. Dissociating fast and slow learning dynamics from the *rich* and *lazy* learning regimes (section 4).
- We analyse how weights dynamically align to task-relevant structure over the course of learning, going beyond prior work that has assumed initial alignment (section 5).
- We provide exact solutions to continual learning dynamics, reversal learning dynamics and to the dynamics of learning and revising structured representations (section 6).

1.2. Related work

Our work builds on analyses of deep linear networks (Baldi and Hornik 1989, Fukumizu 1998, Saxe *et al* 2014, 2019, Lampinen and Ganguli 2018, Arora *et al* 2018a, Tarmoun *et al* 2021, Atanasov *et al* 2022), which have shown that this simple model nevertheless has intricate fixed point structure and nonlinear learning dynamics reminiscent of phenomena seen in nonlinear networks. A variety of works has analysed convergence (Arora *et al* 2018b, Du and Hu 2019), generalisation (Lampinen and Ganguli 2018, Poggio *et al* 2018, Huh 2020), and the implicit bias of gradient descent (Gunasekar *et al* 2018, Ji and Telgarsky 2018, Laurent and Brecht 2018, Arora *et al* 2019a). These works mostly considers the *tabula rasa* case of small initial random weights, for which exact solutions are known (Saxe *et al* 2014). By contrast our formalism describes dynamics from a much larger class of initial conditions and can describe alignment dynamics that do not occur in the *tabula rasa* setting. Most directly, our results build from the matrix Riccati formulation proposed by Fukumizu (1998). Connecting this formulation and matrix factorisation problems yields a better characterisation of the convergence rate (Tarmoun *et al* 2021). We extend and refine the matrix Riccati result to obtain the dynamics of over- and under-complete networks; to obtain numerically stable forms of the matrix equations; and to more explicitly reveal the impact of initialisation.

A line of theoretical research has considered online learning dynamics in teacher-student settings (Biehl and Schwarze 1995, Saad and Solla 1995, Goldt *et al* 2019), deriving ordinary differential equations for the average learning dynamics even in nonlinear networks. However, solving these equations requires numerical integration. By contrast, our approach provides explicit analytical solutions for the more restricted case of deep linear networks.

Other approaches for analysing deep network dynamics include the NTK (Jacot *et al* 2018, Lee *et al* 2019, Arora *et al* 2019b) and the mean field approach (Mei *et al* 2018, Rotskoff and Vanden-Eijnden 2018, Sirignano and Spiliopoulos 2020). While the former can describe nonlinear networks but not the learning dynamics of hidden representations, the later yields a description of representation learning dynamics in wide networks in terms of a partial differential equation. Our work is similar in seeking a subset of more tractable models that are amenable to analysis, but we

focus on the impact of initialisation on representation learning dynamics and explicit solutions.

A large body of work has investigated the effect of different random initialisations on learning in deep networks. The role of initialisation in the vanishing gradient problem and proposals for better initialisation schemes have been illuminated by several works drawing on the central limit theorem (Glorot and Bengio 2010, Saxe *et al* 2014, He *et al* 2015, Pennington *et al* 2017, Xiao *et al* 2018), reviewed in Carleo *et al* (2019), Arora *et al* (2020), Bahri *et al* (2020). These approaches typically guarantee that gradients do not vanish at the start of learning, but do not analytically describe the resulting learning trajectories. Influential work has shown that network initialisation variance mediates a transition from *rich* representation learning to *lazy* NTK dynamics (Chizat *et al* 2019), which we analyse in our framework.

2. Preliminaries and setting

Consider a supervised learning task in which input vectors $\mathbf{x}_n \in \mathbb{R}^{N_i}$ from a set of P training pairs $\{(\mathbf{x}_n, \mathbf{y}_n)\}_{n=1\dots P}$ have to be associated with their target output vectors $\mathbf{y}_n \in \mathbb{R}^{N_o}$. We learn this task with a two-layer linear network model (figure 1(A)) that produces the output prediction

$$\hat{\mathbf{y}}_n = \mathbf{W}_2 \mathbf{W}_1 \mathbf{x}_n, \quad (1)$$

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 N_h is the number of hidden units. The network's weights are optimised using full batch gradient descent with learning rate η (or respectively time constant $\tau = \frac{1}{\eta}$) on the mean squared error loss

$$\mathcal{L}(\hat{\mathbf{y}}, \mathbf{y}) = \frac{1}{2} \langle \|\hat{\mathbf{y}} - \mathbf{y}\|^2 \rangle, \quad (2)$$

where $\langle \cdot \rangle$ denotes the average over the dataset. The input and input-output correlation matrices of the dataset are

$$\tilde{\Sigma}^{xx} = \frac{1}{P} \sum_{n=1}^P \mathbf{x}_n \mathbf{x}_n^T \in \mathbb{R}^{N_i \times N_i} \quad \text{and} \quad \tilde{\Sigma}^{yx} = \frac{1}{P} \sum_{n=1}^P \mathbf{y}_n \mathbf{x}_n^T \in \mathbb{R}^{N_o \times N_i}. \quad (3)$$

Finally, the gradient optimisation starts from an initialisation $\mathbf{W}_2(0), \mathbf{W}_1(0)$. Our goal is to understand the full time trajectory of the network's output and internal representations as a function of this initialisation and the task statistics.

Our starting point is the seminal work of Fukumizu (Fukumizu 1998), which showed that the gradient flow dynamics could be written as a matrix Riccati equation with known solution. In particular, defining

$$\mathbf{Q} = \begin{bmatrix} \mathbf{W}_1^T \\ \mathbf{W}_2 \end{bmatrix} \quad \text{and} \quad \mathbf{F} = \begin{bmatrix} 0 & (\tilde{\Sigma}^{yx})^T \\ \tilde{\Sigma}^{yx} & 0 \end{bmatrix}, \quad (4)$$

the continuous time dynamics of the matrix $\mathbf{Q}\mathbf{Q}^T$ from initial state $\mathbf{Q}(0)$ is

$$\mathbf{Q}\mathbf{Q}^T(t) = e^{\mathbf{F}\frac{t}{\tau}}\mathbf{Q}(0) \left[\mathbf{I} + \frac{1}{2}\mathbf{Q}(0)^T \left(e^{\mathbf{F}\frac{t}{\tau}}\mathbf{F}^{-1}e^{\mathbf{F}\frac{t}{\tau}} - \mathbf{F}^{-1} \right) \mathbf{Q}(0) \right]^{-1} \mathbf{Q}(0)^T e^{\mathbf{F}\frac{t}{\tau}}, \quad (5)$$

if the following four assumptions hold:

Assumption 2.1. The dimensions of the input and target vectors are identical, that is $N_i = N_o$.

Assumption 2.2. The input data is whitened, that is $\tilde{\Sigma}^{xx} = \mathbf{I}$.

Assumption 2.3. The network's weight matrices are zero-balanced at the beginning of training, that is $\mathbf{W}_1(0)\mathbf{W}_1(0)^T = \mathbf{W}_2(0)^T\mathbf{W}_2(0)$. If this condition holds at initialisation, it will persist throughout training (Saxe *et al* 2014, Arora *et al* 2018a).

Assumption 2.4. The input-output correlation of the task and the initial state of the network function have full rank, that is $\text{rank}(\tilde{\Sigma}^{xy}) = \text{rank}(\mathbf{W}_2(0)\mathbf{W}_1(0)) = N_i = N_o$. This implies that the network is not bottlenecked, i.e. $N_h \geq \min(N_i, N_o)$.

For completeness, we include a derivation of this solution in appendix A.

Rather than tracking the weights' dynamics directly, this approach tracks several key statistics collected in the matrix

$$\mathbf{Q}\mathbf{Q}^T = \begin{bmatrix} \mathbf{W}_1^T\mathbf{W}_1(t) & \mathbf{W}_1^T\mathbf{W}_2^T(t) \\ \mathbf{W}_2\mathbf{W}_1(t) & \mathbf{W}_2\mathbf{W}_2^T(t) \end{bmatrix}, \quad (6)$$

which can be separated into four quadrants with intuitive meaning: the off-diagonal blocks contain the network function

$$\hat{\mathbf{Y}}(t) = \mathbf{W}_2\mathbf{W}_1(t)\mathbf{X}, \quad (7)$$

while the on-diagonal blocks contain the correlation structure of the weight matrices. These permit calculation of the temporal evolution of the network's internal representations including the task-relevant representational similarity matrices (RSMs) (Kriegeskorte *et al* 2008), i.e. the kernel matrix $\phi(x)^T\phi(x')$, of the neural representations in the hidden layer

$$\text{RSM}_I = \mathbf{X}^T\mathbf{W}_1^T\mathbf{W}_1(t)\mathbf{X}, \quad \text{RSM}_O = \mathbf{Y}^T(\mathbf{W}_2\mathbf{W}_2^T(t))^+ \mathbf{Y}, \quad (8)$$

where $+$ denotes the pseudoinverse; and the network's finite-width NTK (Jacot *et al* 2018, Lee *et al* 2019, Arora *et al* 2019b)

$$\text{NTK} = \mathbf{I}_{N_o} \otimes \mathbf{X}^T\mathbf{W}_1^T\mathbf{W}_1(t)\mathbf{X} + \mathbf{W}_2\mathbf{W}_2^T(t) \otimes \mathbf{X}^T\mathbf{X}, \quad (9)$$

where \mathbf{I} is the identity matrix and \otimes is the Kronecker product. For a derivation of these quantities see appendix B. Hence, the solution in equation (5) describes important aspects of network behaviour.

However, in this form, the solution has several limitations. First, it relies on general matrix exponentials and inverses, which are a barrier to explicit understanding. Second, when evaluated numerically, it is often unstable. And third, the equation is only valid for equal input and output dimensions. In the following section we address these limitations. Implementation and simulation. Simulation details are in appendix H. Code to replicate all simulations and plots are available online⁶ under a *GPLv3* license and requires <6 h to execute on a single AMD Ryzen 5950x.

3. Exact learning dynamics with prior knowledge

In this section we derive an exact and numerically stable solution for $\mathbf{Q}\mathbf{Q}^T$ that better reveals the learning dynamics, convergence behaviour and generalisation properties of two-layer linear networks with prior knowledge. Further, we alter the equations to be applicable to equal and unequal input and output dimensions, overcoming assumption 2.1.

To place the solution in a more explicit form, we make use of the compact singular value decomposition. Let the compact singular value decomposition of the initial network function and the input-output correlation of the task be

$$\text{SVD}(\mathbf{W}_2(0)\mathbf{W}_1(0)) = \mathbf{U}\mathbf{S}\mathbf{V}^T \quad \text{and} \quad \text{SVD}(\tilde{\Sigma}^{yx}) = \tilde{\mathbf{U}}\tilde{\mathbf{S}}\tilde{\mathbf{V}}^T. \quad (10)$$

Here, \mathbf{U} and $\tilde{\mathbf{U}} \in \mathbb{R}^{N_o \times N_m}$ denote the left singular vectors, \mathbf{S} and $\tilde{\mathbf{S}} \in \mathbb{R}^{N_m \times N_m}$ the square matrix with ordered, non-zero eigenvalues on its diagonal and \mathbf{V} and $\tilde{\mathbf{V}} \in \mathbb{R}^{N_i \times N_m}$ the corresponding right singular vectors. For unequal input-output dimensions ($N_i \neq N_o$) the right and left singular vectors are therefore not generally square and orthonormal. Accordingly, for the case $N_i > N_o$, we define $\tilde{\mathbf{U}}_{\perp} \in \mathbb{R}^{N_o \times (N_o - N_i)}$ as a matrix containing orthogonal column vectors that complete the basis, i.e. make $[\tilde{\mathbf{U}} \tilde{\mathbf{U}}_{\perp}]$ orthonormal. Conversely, we define $\tilde{\mathbf{V}}_{\perp} \in \mathbb{R}^{N_i \times (N_i - N_o)}$ for the case of $N_i > N_o$.

Assumption 3.1. Define $\mathbf{B} = \mathbf{U}^T\tilde{\mathbf{U}} + \mathbf{V}^T\tilde{\mathbf{V}}$ and $\mathbf{C} = \mathbf{U}^T\tilde{\mathbf{U}} - \mathbf{V}^T\tilde{\mathbf{V}}$. \mathbf{B} is non-singular.

Theorem 3.1. Under the assumptions of whitened inputs, 2.2, zero-balanced weights 2.3, full rank 2.4, and \mathbf{B} non-singular 3.1, the temporal dynamics of $\mathbf{Q}\mathbf{Q}^T$ are

$$\begin{aligned} \mathbf{Q}\mathbf{Q}^T(t) = & \mathbf{Z} \left[4e^{-\tilde{\mathbf{S}}_{\tau}^t} \mathbf{B}^{-1} \mathbf{S}^{-1} (\mathbf{B}^T)^{-1} e^{-\tilde{\mathbf{S}}_{\tau}^t} + (\mathbf{I} - e^{-2\tilde{\mathbf{S}}_{\tau}^t}) \tilde{\mathbf{S}}^{-1} \right. \\ & - e^{-\tilde{\mathbf{S}}_{\tau}^t} \mathbf{B}^{-1} \mathbf{C} \left(e^{-2\tilde{\mathbf{S}}_{\tau}^t} - \mathbf{I} \right) \tilde{\mathbf{S}}^{-1} \mathbf{C}^T (\mathbf{B}^T)^{-1} e^{-\tilde{\mathbf{S}}_{\tau}^t} \\ & \left. + 4\frac{t}{\tau} e^{-\tilde{\mathbf{S}}_{\tau}^t} \mathbf{B}^{-1} \left(\mathbf{V}^T \tilde{\mathbf{V}}_{\perp} \tilde{\mathbf{V}}_{\perp}^T \mathbf{V} + \mathbf{U}^T \tilde{\mathbf{U}}_{\perp} \tilde{\mathbf{U}}_{\perp}^T \mathbf{U} \right) (\mathbf{B}^T)^{-1} e^{-\tilde{\mathbf{S}}_{\tau}^t} \right]^{-1} \mathbf{Z}^T \end{aligned} \quad (11)$$

⁶ <https://github.com/saxelab/deep-linear-networks-with-prior-knowledge>.

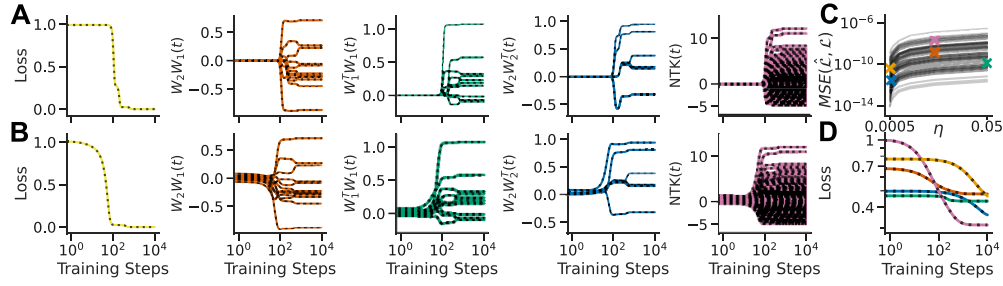


Figure 2. Exact learning dynamics (A) the temporal dynamics of the numerical simulation (coloured lines) of the loss, network function, correlation of input and output weights and the NTK (columns 1–5 respectively) are exactly matched by the analytical solution (black dotted lines) for small initial weight values and (B) large initial weight values. (C) Each line shows the deviation of the analytical loss $\hat{\mathcal{L}}$ from the numerical loss \mathcal{L} for one of $n = 50$ networks with random architecture and training data (details in appendix H) across a range of learning rates $\eta \in [0.05, 0.0005]$. The deviation mutually decreases with the learning rate. (D) Numerical and analytical learning curves for five randomly sampled example networks (coloured x in (C)).

with

$$\mathbf{Z} = \begin{bmatrix} \tilde{\mathbf{V}} \left(\mathbf{I} - e^{-\tilde{\mathbf{S}} \frac{t}{\tau}} \mathbf{C}^T (\mathbf{B}^T)^{-1} e^{-\tilde{\mathbf{S}} \frac{t}{\tau}} \right) + 2\tilde{\mathbf{V}}_{\perp} \tilde{\mathbf{V}}_{\perp}^T \mathbf{V} (\mathbf{B}^T)^{-1} e^{-\tilde{\mathbf{S}} \frac{t}{\tau}} \\ \tilde{\mathbf{U}} \left(\mathbf{I} + e^{-\tilde{\mathbf{S}} \frac{t}{\tau}} \mathbf{C}^T (\mathbf{B}^T)^{-1} e^{-\tilde{\mathbf{S}} \frac{t}{\tau}} \right) + 2\tilde{\mathbf{U}}_{\perp} \tilde{\mathbf{U}}_{\perp}^T \mathbf{U} (\mathbf{B}^T)^{-1} e^{-\tilde{\mathbf{S}} \frac{t}{\tau}} \end{bmatrix}. \quad (12)$$

For a proof of theorem 3.1 please refer to appendix C.

With this solution we can calculate the exact temporal dynamics of the loss, network function, RSMs and NTK (figures 2(A) and (B)). As the solution contains only negative exponentials, it is numerically stable and provides high precision across a wide range of learning rates and network architectures (figures 2(C) and (D)).

We note that a solution for the weights $\mathbf{W}_1(t)$ and $\mathbf{W}_2(t)$, i.e. $\mathbf{Q}(t)$, can be derived up to a time varying orthogonal transformation as demonstrated in appendix C. Further, as time-dependent variables only occur in matrix exponentials of diagonal matrices of negative sign, the network approaches a steady state solution.

Theorem 3.2. *Under the assumptions of theorem 3.1, the network function converges to the global minimum $\tilde{\mathbf{U}}\tilde{\mathbf{S}}\tilde{\mathbf{V}}^T$ and acquires a rich task-specific internal representation, that is $\mathbf{W}_1^T\mathbf{W}_1 = \tilde{\mathbf{V}}\tilde{\mathbf{S}}\tilde{\mathbf{V}}^T$ and $\mathbf{W}_2\mathbf{W}_2^T = \tilde{\mathbf{U}}\tilde{\mathbf{S}}\tilde{\mathbf{U}}^T$.*

The proof of theorem 3.2 is in appendix C. We now turn to several implications of these results.

4. Rich and lazy learning regimes and generalisation

Recent results have shown that large deep networks can operate in qualitatively distinct regimes that depend on their weight initialisations (Chizat *et al* 2019, Flesch *et al* 2022),

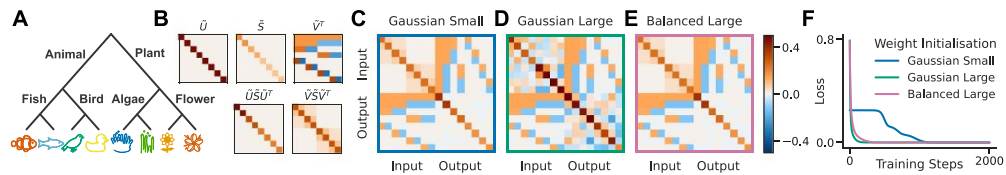


Figure 3. Rich and lazy learning. (A) Semantic learning task, (B) SVD of the input-output correlation of the task (top) and the respective RSMs (bottom). Rows and columns in the SVD and RSMs are identically ordered as the order of items in the hierarchical tree. (C) Final $\mathbf{Q}\mathbf{Q}^T$ matrices after training converged when initialised from random small weights, (D) random large weights (note how the upper left and lower right quadrant differ from the task’s RSMs) and (E) large zero-balanced weights. (F) Learning curves for the three different initialisations as in (C) (green), (D) (pink) and (E) (blue). While both large weight initialisations lead to fast exponential learning curves, the small weight initialisation leads to a slow step-like decay of the loss.

the so called *rich* and *lazy* regimes. In the *rich* regime, learning dynamics can be highly nonlinear and lead to task-specific solutions thought to lead to favourable generalisation properties (Chizat *et al* 2019, Saxe *et al* 2019, Flesch *et al* 2022). By contrast, the *lazy* regime exhibits simple exponential learning dynamics and exploits high-dimensional nonlinear projections of the data produced by the initial random weights, leading to task-agnostic representations that attain zero training error but possibly lower generalisation performance (Jacot *et al* 2018, Lee *et al* 2019, Arora *et al* 2019b). Traditionally, the *rich* and *lazy* learning regimes have been respectively linked to low and high variance initial weights (relative to the network layer size).

To illustrate these phenomena, we consider a semantic learning task in which a set of living things have to be linked to their position in a hierarchical structure (figure 3(A)) (Saxe *et al* 2014). The representational similarity of the input of the task ($\tilde{\mathbf{V}}\tilde{\mathbf{S}}\tilde{\mathbf{V}}^T$) reveals its inherent structure (figure 3(B)). For example, the representations of the two fishes are most similar to each other, less similar to birds and least similar to plants. Likewise, the representational similarity of the task’s target values ($\tilde{\mathbf{U}}\tilde{\mathbf{S}}\tilde{\mathbf{U}}^T$) reveals the primary groups among which items are organised. As a consequence, one can for example predict from an object being a fish that it is an animal and from an object being a plant that it is not a bird. Reflecting these structural relationships in internal representations can allow the *rich* regime to generalise in ways the *lazy* regime cannot. Crucially, $\mathbf{Q}\mathbf{Q}^T(t)$ contains the temporal dynamics of the weights’ representational similarity and therefore can be used to study if a network finds a *rich* or *lazy* solution.

When training a two layer network from random small initial weights, the weights’ input and output RSM (figure 3(C), upper left and lower right quadrant) are identical to the task’s structure at convergence. However, when training from large initial weights, the RSM reveals that the network has converged to a *lazy* solution (figure 3(D)). We emphasise that the network function in both cases is identical (figures 3(C) and (D), lower left quadrant). And while their final loss is identical too, their learning dynamics

evolve slow and step-wise in the case of small initial weights and fast and exponentially in the case of large initial weights (figure 3(F)), as predicted by previous work (Chizat *et al* 2019).

However, from theorem 3.2 it directly follows that our setup is guaranteed to find a *rich* solution in which the weights' RSM is identical to the task's RSM, i.e. $\mathbf{W}_1^T \mathbf{W}_1 = \tilde{\mathbf{V}} \tilde{\mathbf{S}} \tilde{\mathbf{V}}^T$ and $\mathbf{W}_2 \mathbf{W}_2^T = \tilde{\mathbf{U}} \tilde{\mathbf{S}} \tilde{\mathbf{U}}^T$. Therefore, as zero-balanced weights may be large, there exist initial states that converge to *rich* solutions while evolving as rapid exponential learning curves (figures 3(E) and (F)). Crucially, these initialisations are task-agnostic, in the sense that they are independent of the task structure (see Mishkin and Matas 2015). This finding applies to any learning task with well defined input-output correlation. For additional simulations see appendix D. Hence our equation can describe the change in dynamics from step-like to exponential with increasing weight scale, and separate this dynamical phenomenon from the structure of internal representations.

5. Decoupling dynamics

The learning dynamics of deep linear networks depend on the exact initial values of the synaptic weights. Previous solutions studied learning dynamics under the assumption that initial network weights are 'decoupled', such that the initial state of the network and the task share the same singular vectors, i.e. that $\mathbf{U} = \tilde{\mathbf{U}}$ and $\mathbf{V} = \tilde{\mathbf{V}}$ (Saxe *et al* 2014). Intuitively, this assumption means that there is no cross-coupling between different singular modes, such that each evolves independently. However, this assumption is violated in most real-world scenarios. As a consequence, most prior work has relied on the empirical observation that learning from *tabula rasa* small initial weights occurs in two phases: First, the network's input-output map rapidly decouples; then subsequently, independent singular modes are learned in this decoupled regime. Because this decoupling process is fast when training begins from small initial weights, the learning dynamics are still approximately described by the temporal learning dynamics of the singular values assuming decoupling from the start. This dynamic has been called a *silent alignment* process (Atanasov *et al* 2022). Here we leverage our matrix Riccati approach to analytically study the dynamics of this decoupling process. We begin by deriving an alternate form of the exact solution that eases the analysis.

Theorem 5.1. *Let the weight matrices of a two layer linear network be initialised by $\mathbf{W}_1 = \mathbf{A}(0) \tilde{\mathbf{V}}^T$ and $\mathbf{W}_2 = \tilde{\mathbf{U}} \mathbf{A}(0)^T$, where $\mathbf{A}(0) \in \mathbb{R}^{N_i \times N_i}$ is an arbitrary, invertible matrix. Then, under the assumptions of equal input-output dimensions 2.1, whitened inputs 2.2, zero-balanced weights 2.3 and full rank 2.4, the temporal dynamics of $\mathbf{Q}\mathbf{Q}^T$ are fully determined by*

$$\mathbf{A}^T \mathbf{A}(t) = \left[e^{-\tilde{\mathbf{S}} \frac{t}{\tau}} \left(\mathbf{A}(0)^T \mathbf{A}(0) \right)^{-1} e^{-\tilde{\mathbf{S}} \frac{t}{\tau}} + \left(\mathbf{I} - e^{-2\tilde{\mathbf{S}} \frac{t}{\tau}} \right) \tilde{\mathbf{S}}^{-1} \right]^{-1}. \quad (13)$$

Exact learning dynamics of deep linear networks with prior knowledge

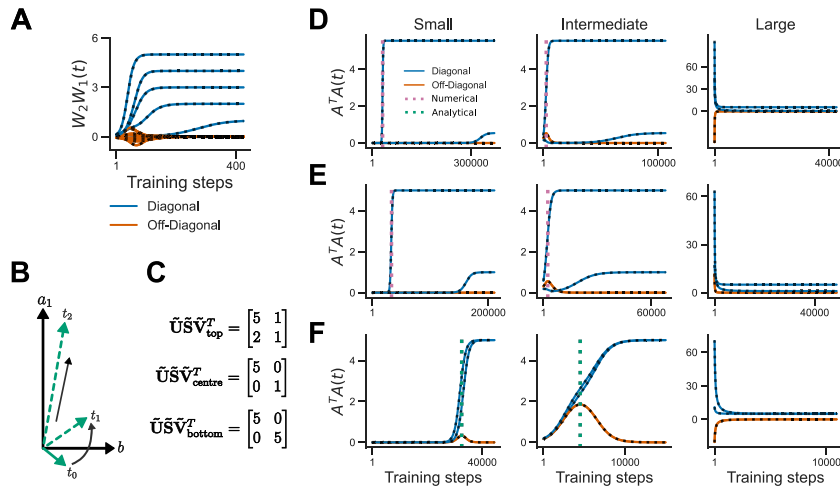


Figure 4. Decoupling dynamics. (A) Analytical (black dotted lines) and numerical (solid lines) of the temporal dynamics of the on- and off-diagonal elements of $\mathbf{A}^T \mathbf{A}$ in blue and red, respectively. (B) Schematic representation of the decoupling process. (C) Three target matrices with dense, unequal diagonal, and equal diagonal structure. (D) and (E) Decoupling dynamics for the top (D), middle (E), and bottom (F) tasks depicted in panel (C). Row F contains analytical predictions for the time of the peak of the off-diagonal (dashed green). The network is initialised as defined in appendix E with small, intermediate and large variance.

For a proof of theorem 5.1, please refer to appendix E. We remark that this form is less general than that in theorem 3.1, and in particular implies $\mathbf{U}\mathbf{V} = \tilde{\mathbf{U}}\tilde{\mathbf{V}}$. Here the matrix $\mathbf{A}^T \mathbf{A}$ represents the dynamics directly in the SVD basis of the task. Off-diagonal elements represent counterproductive coupling between different singular modes (for instance, $[\mathbf{A}^T \mathbf{A}]_{21}$ is the strength of connection from input singular vector 1 to output singular vector 2, which must approach zero to perform the task perfectly), while on-diagonal elements represent the coupling within the same mode (for instance, $[\mathbf{A}^T \mathbf{A}]_{11}$ is the strength of connection from input singular vector 1 to output singular vector 1, which must approach the associated task singular value to perform the task perfectly). Hence the decoupling process can be studied by examining the dynamics by which $\mathbf{A}^T \mathbf{A}$ becomes approximately diagonal.

The outer inverse in equation (13) renders it difficult to study high dimensional networks analytically. Therefore, we focus on small networks with input and output dimension $N_i = 2$ and $N_o = 2$, for which a lengthy but explicit analytical solution is given in appendix E. In this setting, the structure of the weight initialisation and task are encoded in the matrices

$$\mathbf{A}(0)^T \mathbf{A}(0) = \begin{bmatrix} a_1(0) & b(0) \\ b(0) & a_2(0) \end{bmatrix} \quad \text{and} \quad \tilde{\mathbf{S}} = \begin{bmatrix} s_1 & 0 \\ 0 & s_2 \end{bmatrix}, \quad (14)$$

where the parameters $a_1(0)$ and $a_2(0)$ represent the component of the initialisation that is aligned with the task, and $b(0)$ represents cross-coupling, such that taking $b(0) = 0$

recovers previously known and more restricted solutions for the decoupled case (Saxe *et al* 2014). We use this setting to demonstrate two features of the learning dynamics.

Decoupling dynamics. First, we track decoupling by considering the dynamics of the off-diagonal element $b(t)$ (figures 4(D)–(F) red lines). At convergence, the off-diagonal element shrinks to zero as shown in appendix E. However, strikingly, $b(t)$ can exhibit non-monotonic trajectories with transient peaks or valleys partway through the learning process. In particular, in appendix E we derive the time of the peak magnitude as $t_{\text{peak}} = \frac{\tau}{4s} \ln \frac{s(s-a_1-a_2)}{a_1a_2-b(0)^2}$ (figure 4(F) green dotted line), which coincides approximately with the time at which the on-diagonal element is half learned. If initialised from small random weights, the off-diagonal remains near-zero throughout learning, reminiscent of the silent alignment effect (Atanasov *et al* 2022). For large initialisations, no peak is observed and the dynamics are exponential. At intermediate initialisations, the maximum of the off-diagonal is reached before the singular mode is fully learned (appendix E). Intuitively, a particular input singular vector can initially project appreciably onto the wrong output singular vector, corresponding to initial misalignment. This is only revealed when this link is amplified, at which point corrective dynamics remove the counter-productive coupling, as schematised in figure 4(B). We report further measurements of decoupling in appendix E.

Effect of initialisation variance. Next, we revisit the impact of initialisation scale for the on-diagonal dynamics. As shown in figures 4(D)–(F), as the initialisation variance grows the learning dynamics change from sigmoidal to exponential, possibly displaying more complex behaviour at intermediate variance (appendix E). In this simple setting we can analyse this transition in detail. Taking $s_1 = s_2 = s$ as in figure 4(F) and $|a_1(0)|, |a_2(0)|, |b(0)| \ll 1$, we recover a sigmoidal trajectory,

$$a_1(t) = \frac{sa_1(0)}{e^{\frac{-2st}{\tau}} [s - a_1(0) - a_2(0)] + a_1(0) + a_2(0)}, \quad (15)$$

while for $|a_1(0)|, |a_2(0)|, |b(0)| \gg 0$ the dynamics of the on-diagonal element a_1 is close to exponential (figures 4(D)–(F) left and right columns). We examine larger networks in appendix E.

6. Applications

The solutions derived in sections 3 and 5 provide tools to examine the impact of prior knowledge on dynamics in deep linear networks. So far we have traced general features of the behaviour of these solutions. In this section, we use this toolkit to develop accounts of several specific phenomena.

Continual Learning. Continual learning (see Parisi *et al* 2019 for a review) and the pathology of catastrophic forgetting have long been a challenge for neural network models (McCloskey and Cohen 1989, Ratcliff 1990, French 1999). A variety of theoretical work has investigated aspects of continual learning (Tripuraneni *et al* 2020, Asanuma

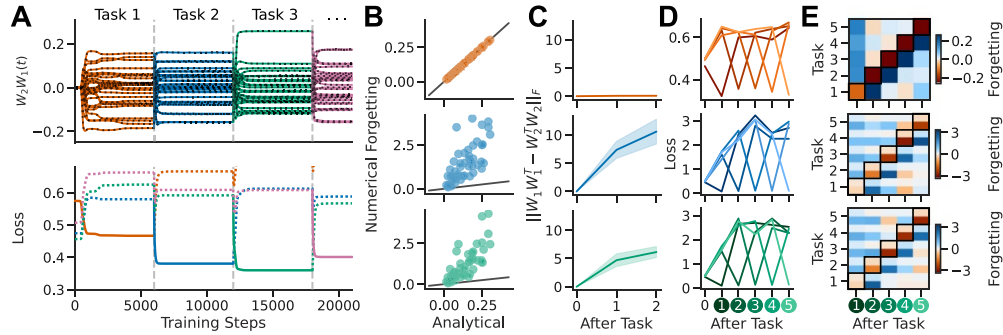


Figure 5. Continual learning. (A) Top: network training from small zero-balanced weights on a sequence of tasks (coloured lines show simulation and black dotted lines analytical results). Bottom: evaluation loss for tasks of the sequence (dotted) while training on the current task (solid). As the network function is optimised on the current task, the loss of other tasks increases. (B) Comparison of the numerical and analytical amount of catastrophic forgetting on a first task after training on a second task for $n = 50$ linear (red), tanh (blue) and ReLU (green) networks. (C) Weight alignment before and after training on a sequence of two tasks for $n = 50$ networks in linear (red), tanh (blue) and ReLU (green) networks. Shaded area shows \pm std. (D) Evaluation loss for each of 5 tasks during training a linear (red), tanh (blue) and ReLU (green) network. (E) Same data as in (D) but evaluated as relative change (i.e. amount of catastrophic forgetting). The top half of each square shows the pre-computed analytical amount of forgetting and the bottom half the numerical value.

et al 2021, Doan *et al* 2021, Lee *et al* 2021, Shachaf *et al* 2021). In this setting, starting from an initial set of weights, a network is trained on a sequence of tasks with respective input-output correlations $\mathcal{T}_1 = \tilde{\Sigma}_1^{yx}, \mathcal{T}_2 = \tilde{\Sigma}_2^{yx}, \mathcal{T}_3 = \tilde{\Sigma}_3^{yx}, \dots$. As shown in figure 5(A), our dynamics immediately enable exact solutions for the full continual learning process, whereby the final state after training on one task becomes the initial network state for the next task. These solutions thus reveal the exact time course of forgetting for arbitrary sequences of tasks.

Training on later tasks can overwrite previously learned knowledge, a phenomenon known as catastrophic forgetting (McCloskey and Cohen 1989, Ratcliff 1990, French 1999). From theorem 3.2 it follows that from any arbitrary zero-balanced initialisation 2.3, the network converges to the global optimum such that the initialisation is completely overwritten and forgetting is truly catastrophic. In particular, the loss of any other task \mathcal{T}_i after training to convergence on task \mathcal{T}_j is $\mathcal{L}_i(\mathcal{T}_j) = 1/2 \|\tilde{\Sigma}_j^{yx} - \tilde{\Sigma}_i^{yx}\|_F^2 + c$, where c is a constant that only depends on training data of task \mathcal{T}_i (appendix F). As a consequence, the amount of forgetting, i.e. the relative change of loss, is fully determined by the similarity structure of the tasks and thus can be fully determined for a sequence of tasks before the onset of training (figures 5(B) and (E), appendix F). For example, the amount of catastrophic forgetting in task \mathcal{T}_a , when training on task \mathcal{T}_c after having trained the network on task \mathcal{T}_b is $\mathcal{L}_a(\mathcal{T}_c) - \mathcal{L}_a(\mathcal{T}_b)$. As expected, our results depend on our linear setting and tanh or ReLU nonlinearities can show different behaviour, typically increasing the amount of forgetting (figures 5(B), (D) and (E)). Further, in

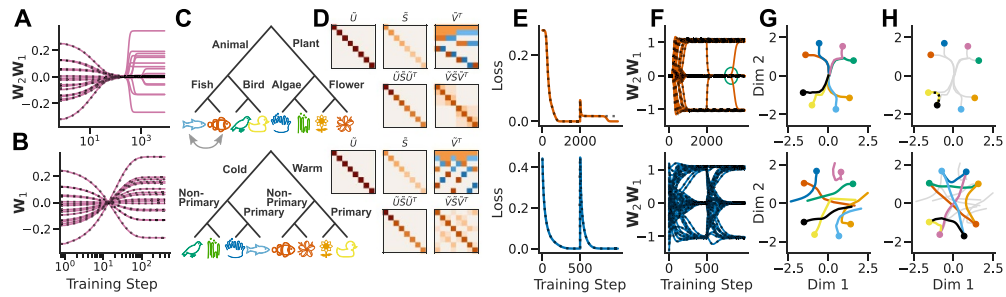


Figure 6. Reversal learning and revising structured knowledge. Scale of x -axis varies in top and bottom rows. (A) Analytical (black dotted) and numerical (solid) learning dynamics of a reversal learning task. The analytical solution gets stuck on a saddle point, whereas the numerical simulation escapes the saddle point and converges to the target. (B) In a shallow network, training on the same task as in A converges analytically (black dotted) and numerically (solid). (C) Semantic learning tasks. Revised living kingdom (top) and colour hierarchy (bottom). (D) SVD of the input-output correlation of the tasks and respective RSMs. (E) Analytical (black dotted) and simulation (solid) loss and (F) learning dynamics of first training on the living kingdom (figure 3(A)) and subsequently on the respective task in (C). The analytical solution fails for the revised animal kingdom as it gets stuck in a saddle point, while the simulation escapes the saddle (top, green circle). Initial training on the living kingdom task from large initial weights and subsequent training on the colour hierarchy have similar convergence times (bottom) (G) multidimensional scaling (MDS) of the network function for initial training on the living kingdom task from small (top) and large initial weights (bottom). Note how despite the seemingly chaotic learning dynamics when starting from large initial weights, both simulations learn the same representation. (H) MDS of subsequent training on the respective task in (C).

nonlinear networks, weights become rapidly unbalanced and forgetting values that are calculated before the onset of training do not predict the actual outcome (figures 5(B)–(E)). In summary, our results link exact learning dynamics with catastrophic forgetting and thus provide an analytical tool to study the mechanisms and potential counter measures underlying catastrophic forgetting.

Reversal learning. During reversal learning, pre-existing knowledge has to be relearned, overcoming a previously learned relationship between inputs and outputs. For example, reversal learning occurs when items of a class are mislabelled and later corrected. We show analytically, that reversal learning in fact does not succeed in deep linear networks (appendix G). The pre-existing knowledge lies exactly on the separatrix of a saddle point causing the learning dynamics to converge to zero (figure 6(A)). In contrast, the learning still succeeds numerically, as any noise will perturb the dynamics off the saddle point, allowing learning to proceed (figure 6(A)). However, the dynamics still slow in the vicinity of the saddle point, providing a theoretical explanation for catastrophic slowing in deep linear networks (Lee *et al* 2022). We note that the analytical

solution requires an adaptation of theorem 3.1, as \mathbf{B} is generally not invertible in the case of reversal learning (appendix G). Further, as is revealed by the exact learning dynamics (appendix G), shallow networks do succeed without exhibiting catastrophic slowing during reversal learning (figure 6(B)).

Revising structured knowledge. Knowledge is often organised within an underlying, shared structure, of which many can be learned and represented in deep linear networks (Saxe *et al* 2019). For example, spatial locations can be related to each other using the same cardinal directions, or varying semantic knowledge can be organised using the same hierarchical tree. Here, we investigate if deep linear networks benefit from shared underlying structure. To this end, a network is first trained on the three-level hierarchical tree of section 4 (eight items of the living kingdom, each with a set of eight associated features), and subsequently trained on a revised version of the hierarchy. The revised task varies the relation of inputs and outputs while keeping the same underlying tree structure. If the revision involves swapping two neighbouring nodes on any level of the hierarchy, e.g. the identity of the two fish on the lowest level of the hierarchy (figure 6(C), top), the task is identical to reversal learning, leading to catastrophically slowed dynamics (figures 6(E) and (F), top). When training the network on a new hierarchical tree with identical items but a new set of features, like a colour hierarchy (figure 6(C), bottom), there is no speed advantage in comparison to a random initialisation with similar initial variance (figures 6(E) and (F), bottom). Importantly, from theorem 3.2 it follows, that the learning process can be sped up significantly by initialising from large zero-balanced weights, while converging to a global minimum with identical generalisation properties as when training from small weights (figures 6(G) and (H)). In summary, having incorporated structured knowledge before revision does not speed up or even slows down learning in comparison to learning from random zero-balanced weights. Notably, that is despite the tasks' structure being almost identical (figures 3(B) and 6(D)).

7. Discussion

We derive exact solutions to the dynamics of learning with rich prior knowledge in a tractable model class: deep linear networks. While our results broaden the class of two-layer linear network problems that can be described analytically, they remain limited and rely on a set of assumptions (2.1)–(2.4). In particular, weakening the requirement that the input covariance be white and the weights be zero-balanced would enable analysis of the impact of initialisation on internal representations. Nevertheless, these solutions reveal several insights into network behaviour. We show that there exists a large set of initial values, namely zero-balanced weights 2.3, which lead to task-specific representations; and that large initialisations lead to exponential rather than sigmoidal learning curves. We hope our results provide a mathematical toolkit that illuminates the complex impact of prior knowledge on deep learning dynamics.

Acknowledgment

L B was supported by the Woodward Scholarship awarded by Wadham College, Oxford and the Medical Research Council [MR/N013468/1]. C D and A S were supported by the Gatsby Charitable Foundation (GAT3755). Further, A S was supported by a Sir Henry Dale Fellowship from the Wellcome Trust and Royal Society (216386/Z/19/Z) and the Sainsbury Wellcome Centre Core Grant (219627/Z/19/Z). A S is a CIFAR Azrieli Global Scholar in the Learning in Machines & Brains program. J F was supported by the Howard Hughes Medical Institute.

Appendix A. Fukumizu approach

For completeness, we reproduce the derivation from Fukumizu (1998) of equation (5). We consider the learning setting describe in section 2. Under the assumptions of equal input-output dimensions 2.1, whitened inputs 2.2 and zero-balanced weights 2.3, the weights dynamics yield

$$\tau \frac{d}{dt} \mathbf{W}_1 = \mathbf{W}_2^T \left(\tilde{\Sigma}^{yx} - \mathbf{W}_2 \mathbf{W}_1 \tilde{\Sigma}^{xx} \right), \quad (16)$$

$$\tau \frac{d}{dt} \mathbf{W}_2 = \left(\tilde{\Sigma}^{yx} - \mathbf{W}_2 \mathbf{W}_1 \tilde{\Sigma}^{xx} \right) \mathbf{W}_1^T. \quad (17)$$

Under the assumption of whitened inputs 2.2, the dynamics simplify to

$$\tau \frac{d}{dt} \mathbf{W}_1 = \mathbf{W}_2^T \left(\tilde{\Sigma}^{yx} - \mathbf{W}_2 \mathbf{W}_1 \right), \quad (18)$$

$$\tau \frac{d}{dt} \mathbf{W}_2 = \left(\tilde{\Sigma}^{yx} - \mathbf{W}_2 \mathbf{W}_1 \right) \mathbf{W}_1^T. \quad (19)$$

We introduce the variables

$$\mathbf{Q} = \begin{bmatrix} \mathbf{W}_1^T \\ \mathbf{W}_2 \end{bmatrix} \text{ and } \mathbf{Q}\mathbf{Q}^T = \begin{bmatrix} \mathbf{W}_1^T \mathbf{W}_1 & \mathbf{W}_1^T \mathbf{W}_2^T \\ \mathbf{W}_2 \mathbf{W}_1 & \mathbf{W}_2 \mathbf{W}_2^T \end{bmatrix}. \quad (20)$$

We compute the time derivative

$$\tau \frac{d}{dt} (\mathbf{Q}\mathbf{Q}^T) = \tau \begin{bmatrix} \frac{d\mathbf{W}_1^T}{dt} \mathbf{W}_1 + \mathbf{W}_1^T \frac{d\mathbf{W}_1}{dt} & \frac{d\mathbf{W}_1^T}{dt} \mathbf{W}_2^T + \mathbf{W}_1^T \frac{d\mathbf{W}_2^T}{dt} \\ \frac{d\mathbf{W}_2}{dt} \mathbf{W}_1 + \mathbf{W}_2 \frac{d\mathbf{W}_1}{dt} & \frac{d\mathbf{W}_2}{dt} \mathbf{W}_2^T + \mathbf{W}_2 \frac{d\mathbf{W}_2^T}{dt} \end{bmatrix}. \quad (21)$$

Using equations (18) and (19) we compute the four quadrant separately giving

$$\tau \left(\frac{d\mathbf{W}_1^T}{dt} \mathbf{W}_1 + \mathbf{W}_1^T \frac{d\mathbf{W}_1}{dt} \right) \quad (22)$$

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

$$= \left(\tilde{\Sigma}^{yx}\right)^T \mathbf{W}_2 \mathbf{W}_1 + \mathbf{W}_1^T \mathbf{W}_2^T \tilde{\Sigma}^{yx} - \mathbf{W}_1^T \mathbf{W}_2^T \mathbf{W}_2 \mathbf{W}_1 - (\mathbf{W}_2 \mathbf{W}_1)^T \mathbf{W}_2 \mathbf{W}_1 \quad (24)$$

$$= \left(\tilde{\Sigma}^{yx}\right)^T \mathbf{W}_2 \mathbf{W}_1 + \mathbf{W}_1^T \mathbf{W}_2^T \tilde{\Sigma}^{yx} - \mathbf{W}_1^T \mathbf{W}_2^T \mathbf{W}_2 \mathbf{W}_1 - \mathbf{W}_1^T \mathbf{W}_1 \mathbf{W}_1^T \mathbf{W}_1, \quad (25)$$

$$\tau \left(\frac{d\mathbf{W}_1^T}{dt} \mathbf{W}_2^T + \mathbf{W}_1^T \frac{d\mathbf{W}_2^T}{dt} \right) \quad (26)$$

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

$$= \left(\tilde{\Sigma}^{yx}\right)^T \mathbf{W}_2 \mathbf{W}_2^T + \mathbf{W}_1^T \mathbf{W}_1 \left(\tilde{\Sigma}^{yx}\right)^T - \mathbf{W}_1^T \mathbf{W}_1 (\mathbf{W}_2 \mathbf{W}_1)^T - (\mathbf{W}_2 \mathbf{W}_1)^T \mathbf{W}_2 \mathbf{W}_2^T, \quad (28)$$

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

$$\tau \left(\frac{d\mathbf{W}_2}{dt} \mathbf{W}_1 + \mathbf{W}_2 \frac{d\mathbf{W}_1}{dt} \right) \quad (30)$$

$$= \left(\tilde{\Sigma}^{yx} - \mathbf{W}_2 \mathbf{W}_1\right) \mathbf{W}_1^T \mathbf{W}_1 + \mathbf{W}_2 \mathbf{W}_2^T \left(\tilde{\Sigma}^{yx} - \mathbf{W}_2 \mathbf{W}_1\right) \quad (31)$$

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

$$\tau \left(\frac{d\mathbf{W}_2}{dt} \mathbf{W}_2^T + \mathbf{W}_2 \frac{d\mathbf{W}_2^T}{dt} \right) \quad (33)$$

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

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

$$= \tilde{\Sigma}^{yx} \mathbf{W}_1^T \mathbf{W}_2^T + \mathbf{W}_2 \mathbf{W}_1 \left(\tilde{\Sigma}^{yx}\right)^T - \mathbf{W}_2 \mathbf{W}_1 \mathbf{W}_1^T \mathbf{W}_2^T - \mathbf{W}_2 \mathbf{W}_1 \mathbf{W}_1^T \mathbf{W}_2^T \quad (36)$$

$$= \tilde{\Sigma}^{yx} \mathbf{W}_1^T \mathbf{W}_2^T + \mathbf{W}_2 \mathbf{W}_1 \left(\tilde{\Sigma}^{yx}\right)^T - \mathbf{W}_2 \mathbf{W}_1 \mathbf{W}_1^T \mathbf{W}_2^T - \mathbf{W}_2 \mathbf{W}_2^T \mathbf{W}_2 \mathbf{W}_2^T, \quad (37)$$

where we have used the assumption of zero-balanced weights 2.3 to simplify equations (25) and (37).

Defining

$$\mathbf{F} = \begin{bmatrix} 0 & \left(\tilde{\Sigma}^{yx}\right)^T \\ \tilde{\Sigma}^{yx} & 0 \end{bmatrix}, \quad (38)$$

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

$$\tau \frac{d}{dt} (\mathbf{Q}\mathbf{Q}^T) = \mathbf{F}\mathbf{Q}\mathbf{Q}^T + \mathbf{Q}\mathbf{Q}^T\mathbf{F} - (\mathbf{Q}\mathbf{Q}^T)^2. \tag{39}$$

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

$$\begin{aligned} &\tau \frac{d}{dt} (\mathbf{Q}\mathbf{Q}^T) \\ &= \begin{bmatrix} 0 & (\tilde{\Sigma}^{yx})^T \\ \tilde{\Sigma}^{yx} & 0 \end{bmatrix} \begin{bmatrix} \mathbf{W}_1^T \mathbf{W}_1 & \mathbf{W}_1^T \mathbf{W}_2^T \\ \mathbf{W}_2 \mathbf{W}_1 & \mathbf{W}_2 \mathbf{W}_2^T \end{bmatrix} \\ &\quad + \begin{bmatrix} \mathbf{W}_1^T \mathbf{W}_1 & \mathbf{W}_1^T \mathbf{W}_2^T \\ \mathbf{W}_2 \mathbf{W}_1 & \mathbf{W}_2 \mathbf{W}_2^T \end{bmatrix}^T \begin{bmatrix} 0 & (\tilde{\Sigma}^{yx})^T \\ \tilde{\Sigma}^{yx} & 0 \end{bmatrix} - \begin{bmatrix} \mathbf{W}_1^T \mathbf{W}_1 & \mathbf{W}_1^T \mathbf{W}_2^T \\ \mathbf{W}_2 \mathbf{W}_1 & \mathbf{W}_2 \mathbf{W}_2^T \end{bmatrix}^2 \end{aligned} \tag{40}$$

$$\begin{aligned} &= \begin{bmatrix} 0 & (\tilde{\Sigma}^{yx})^T \\ \tilde{\Sigma}^{yx} & 0 \end{bmatrix} \begin{bmatrix} \mathbf{W}_1^T \mathbf{W}_1 & \mathbf{W}_1^T \mathbf{W}_2^T \\ \mathbf{W}_2 \mathbf{W}_1 & \mathbf{W}_2 \mathbf{W}_2^T \end{bmatrix} \\ &\quad + \begin{bmatrix} \mathbf{W}_1^T \mathbf{W}_1 & \mathbf{W}_1^T \mathbf{W}_2^T \\ \mathbf{W}_2 \mathbf{W}_1 & \mathbf{W}_2 \mathbf{W}_2^T \end{bmatrix} \begin{bmatrix} 0 & (\tilde{\Sigma}^{yx})^T \\ \tilde{\Sigma}^{yx} & 0 \end{bmatrix} \\ &\quad - \begin{bmatrix} \mathbf{W}_1^T \mathbf{W}_1 & \mathbf{W}_1^T \mathbf{W}_2^T \\ \mathbf{W}_2 \mathbf{W}_1 & \mathbf{W}_2 \mathbf{W}_2^T \end{bmatrix} \begin{bmatrix} \mathbf{W}_1^T \mathbf{W}_1 & \mathbf{W}_1^T \mathbf{W}_2^T \\ \mathbf{W}_2 \mathbf{W}_1 & \mathbf{W}_2 \mathbf{W}_2^T \end{bmatrix} \end{aligned} \tag{41}$$

$$\begin{aligned} &= \begin{bmatrix} (\tilde{\Sigma}^{yx})^T \mathbf{W}_2 \mathbf{W}_1 & (\tilde{\Sigma}^{yx})^T \mathbf{W}_2 \mathbf{W}_2^T \\ \tilde{\Sigma}^{yx} \mathbf{W}_1^T \mathbf{W}_1 & \tilde{\Sigma}^{yx} \mathbf{W}_1^T \mathbf{W}_2^T \end{bmatrix} \\ &\quad + \begin{bmatrix} \mathbf{W}_1^T \mathbf{W}_2^T \tilde{\Sigma}^{yx} & \mathbf{W}_1^T \mathbf{W}_1 (\tilde{\Sigma}^{yx})^T \\ \mathbf{W}_2 \mathbf{W}_2^T \tilde{\Sigma}^{yx} & \mathbf{W}_2 \mathbf{W}_1 (\tilde{\Sigma}^{yx})^T \end{bmatrix} \\ &\quad - \begin{bmatrix} \mathbf{W}_1^T \mathbf{W}_1 & \mathbf{W}_1^T \mathbf{W}_2^T \\ \mathbf{W}_2 \mathbf{W}_1 & \mathbf{W}_2 \mathbf{W}_2^T \end{bmatrix} \begin{bmatrix} \mathbf{W}_1^T \mathbf{W}_1 & \mathbf{W}_1^T \mathbf{W}_2^T \\ \mathbf{W}_2 \mathbf{W}_1 & \mathbf{W}_2 \mathbf{W}_2^T \end{bmatrix} \end{aligned} \tag{42}$$

$$\begin{aligned} &= \begin{bmatrix} (\tilde{\Sigma}^{yx})^T \mathbf{W}_2 \mathbf{W}_1 & (\tilde{\Sigma}^{yx})^T \mathbf{W}_2 \mathbf{W}_2^T \\ \tilde{\Sigma}^{yx} \mathbf{W}_1^T \mathbf{W}_1 & \tilde{\Sigma}^{yx} \mathbf{W}_1^T \mathbf{W}_2^T \end{bmatrix} \\ &\quad + \begin{bmatrix} \mathbf{W}_1^T \mathbf{W}_2^T \tilde{\Sigma}^{yx} & \mathbf{W}_1^T \mathbf{W}_1 (\tilde{\Sigma}^{yx})^T \\ \mathbf{W}_2 \mathbf{W}_2^T \tilde{\Sigma}^{yx} & \mathbf{W}_2 \mathbf{W}_1 (\tilde{\Sigma}^{yx})^T \end{bmatrix} \\ &\quad - \begin{bmatrix} \mathbf{W}_1^T \mathbf{W}_1 \mathbf{W}_1^T \mathbf{W}_1 + \mathbf{W}_1^T \mathbf{W}_2^T \mathbf{W}_2 \mathbf{W}_1 & \mathbf{W}_1^T \mathbf{W}_1 \mathbf{W}_1^T \mathbf{W}_2^T + \mathbf{W}_1^T \mathbf{W}_2^T \mathbf{W}_2 \mathbf{W}_2^T \\ \mathbf{W}_2 \mathbf{W}_1 \mathbf{W}_1^T \mathbf{W}_1 + \mathbf{W}_2 \mathbf{W}_2^T \mathbf{W}_2 \mathbf{W}_1 & \mathbf{W}_2 \mathbf{W}_1 \mathbf{W}_1^T \mathbf{W}_2^T + \mathbf{W}_2 \mathbf{W}_2^T \mathbf{W}_2 \mathbf{W}_2^T \end{bmatrix} \end{aligned} \tag{43}$$

$$= \begin{bmatrix} \left(\tilde{\Sigma}^{yx}\right)^T \mathbf{W}_2 \mathbf{W}_1 + \mathbf{W}_1^T \mathbf{W}_2^T \tilde{\Sigma}^{yx} & \left(\tilde{\Sigma}^{yx}\right)^T \mathbf{W}_2 \mathbf{W}_2^T + \mathbf{W}_1^T \mathbf{W}_1 \left(\tilde{\Sigma}^{yx}\right)^T \\ -\mathbf{W}_1^T \mathbf{W}_2^T \mathbf{W}_2 \mathbf{W}_1 - \mathbf{W}_1^T \mathbf{W}_1 \mathbf{W}_1^T \mathbf{W}_1 & -\mathbf{W}_1^T \mathbf{W}_1 \mathbf{W}_1^T \mathbf{W}_2^T - \mathbf{W}_1^T \mathbf{W}_2^T \mathbf{W}_2 \mathbf{W}_1^T \\ \tilde{\Sigma}^{yx} \mathbf{W}_1^T \mathbf{W}_1 + \mathbf{W}_2 \mathbf{W}_2^T \tilde{\Sigma}^{yx} & \tilde{\Sigma}^{yx} \mathbf{W}_1^T \mathbf{W}_2^T + \mathbf{W}_2 \mathbf{W}_1 \left(\tilde{\Sigma}^{yx}\right)^T \\ -\mathbf{W}_2 \mathbf{W}_2^T \mathbf{W}_2 \mathbf{W}_1 - \mathbf{W}_2 \mathbf{W}_1 \mathbf{W}_1^T \mathbf{W}_1 & -\mathbf{W}_2 \mathbf{W}_1 \mathbf{W}_1^T \mathbf{W}_2^T - \mathbf{W}_2 \mathbf{W}_2^T \mathbf{W}_2 \mathbf{W}_1^T \end{bmatrix} \quad (44)$$

□

The four quadrant of (44) are equivalent to equations (25), (29), (32) and (37) respectively.

Assuming that $\mathbf{Q}(0)$ is full rank, the continuous differential equation (39) has a unique solution for all $t \geq 0$

$$\mathbf{Q}\mathbf{Q}^T(t) = e^{\mathbf{F}\frac{t}{\tau}} \mathbf{Q}(0) \left[\mathbf{I} + \frac{1}{2} \mathbf{Q}(0)^T \left(e^{\mathbf{F}\frac{t}{\tau}} \mathbf{F}^{-1} e^{\mathbf{F}\frac{t}{\tau}} - \mathbf{F}^{-1} \right) \mathbf{Q}(0) \right]^{-1} \mathbf{Q}(0)^T e^{\mathbf{F}\frac{t}{\tau}}. \quad (45)$$

Appendix B. Network’s internal representations

B.1. Representational similarity analysis

The task-relevant representational similarity matrix (Kriegeskorte *et al* 2008) of the hidden layer, calculated from the inputs $\mathbf{H} = \mathbf{W}_1 \mathbf{X}$ is

$$\text{RSM}_I(t) = \mathbf{H}^T(t) \mathbf{H}(t) \quad (46)$$

$$= (\mathbf{W}_1(t) \mathbf{X})^T \mathbf{W}_1(t) \mathbf{X} \quad (47)$$

$$= \mathbf{X}^T (\mathbf{W}_1^T \mathbf{W}_1)(t) \mathbf{X}. \quad (48)$$

Similarly, the representational similarity matrix of the hidden layer, calculated from the outputs $\tilde{\mathbf{H}} = \mathbf{W}_2^+ Y$, where $+$ denotes the pseudoinverse, is

$$\text{RSM}_O(t) = \tilde{\mathbf{H}}^T(t) \tilde{\mathbf{H}}(t) \quad (49)$$

$$= (\mathbf{W}_2^+(t) Y)^T \mathbf{W}_2^+(t) Y \quad (50)$$

$$= Y^T (\mathbf{W}_2 \mathbf{W}_2^T(t))^+ Y. \quad (51)$$

B.2. Finite-width neural tangent kernel

In the following, we derive the finite-width neural tangent kernel (Jacot *et al* 2018) for a two-layer linear network. Starting with the network function at time t

$$F_t(\mathbf{X}) = \mathbf{W}_2 \mathbf{W}_1 \mathbf{X}, \quad (52)$$

the discrete time gradient descent dynamics of the next time step yields

$$F_{t+1}(\mathbf{X}) = \left(\mathbf{W}_2 - \eta \frac{\partial \mathcal{L}}{\partial \mathbf{W}_2} \right) \left(\mathbf{W}_1 - \eta \frac{\partial \mathcal{L}}{\partial \mathbf{W}_1} \right) \mathbf{X} \quad (53)$$

$$= \mathbf{W}_2 \mathbf{W}_1 \mathbf{X} - \eta \left(\mathbf{W}_2 \frac{\partial \mathcal{L}}{\partial \mathbf{W}_1} + \frac{\partial \mathcal{L}}{\partial \mathbf{W}_2} \mathbf{W}_1 - \eta \frac{\partial \mathcal{L}}{\partial \mathbf{W}_2} \frac{\partial \mathcal{L}}{\partial \mathbf{W}_1} \right) \mathbf{X}. \quad (54)$$

The network function's gradient flow can then be derived as

$$\frac{F_{t+1}(\mathbf{X}) - F_t(\mathbf{X})}{\eta} = - \left(\mathbf{W}_2 \frac{\partial \mathcal{L}}{\partial \mathbf{W}_1} + \frac{\partial \mathcal{L}}{\partial \mathbf{W}_2} \mathbf{W}_1 - \eta \frac{\partial \mathcal{L}}{\partial \mathbf{W}_2} \frac{\partial \mathcal{L}}{\partial \mathbf{W}_1} \right) \mathbf{X} \quad (55)$$

$$[\eta \rightarrow 0] \frac{d}{dt} F(\mathbf{X}) = - \left(\mathbf{W}_2 \frac{\partial \mathcal{L}}{\partial \mathbf{W}_1} + \frac{\partial \mathcal{L}}{\partial \mathbf{W}_2} \mathbf{W}_1 \right) \mathbf{X}. \quad (56)$$

Substituting the partial derivatives

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}_1} = \frac{1}{2} \frac{\partial}{\partial \mathbf{W}_1} \|\mathbf{W}_2 \mathbf{W}_1 \mathbf{X} - \mathbf{Y}\|_F^2 \quad (57)$$

$$= \mathbf{W}_2^T (\mathbf{W}_2 \mathbf{W}_1 \mathbf{X} - \mathbf{Y}) \mathbf{X}^T \quad (58)$$

and

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}_2} = \frac{1}{2} \frac{\partial}{\partial \mathbf{W}_2} \|\mathbf{W}_2 \mathbf{W}_1 \mathbf{X} - \mathbf{Y}\|_F^2 \quad (59)$$

$$= (\mathbf{W}_2 \mathbf{W}_1 \mathbf{X} - \mathbf{Y}) \mathbf{X}^T \mathbf{W}_1^T \quad (60)$$

then yields

$$\frac{d}{dt} F(\mathbf{X}) = -\mathbf{W}_2 \mathbf{W}_2^T (\mathbf{W}_2 \mathbf{W}_1 \mathbf{X} - \mathbf{Y}) \mathbf{X}^T \mathbf{X} - (\mathbf{W}_2 \mathbf{W}_1 \mathbf{X} - \mathbf{Y}) \mathbf{X}^T \mathbf{W}_1^T \mathbf{W}_1 \mathbf{X}. \quad (61)$$

Finally, we introduce the identity matrix \mathbf{I}_{N_0} of size N_0 and apply row-wise vectorisation $\text{vec}_r(F(\mathbf{X})) := f(\mathbf{X})$ and the identity $\text{vec}_r(ABC) = (A \otimes C^T) \text{vec}_r(B)$ to derive the neural tangent kernel

$$\frac{d}{dt} F(\mathbf{X}) = -\mathbf{W}_2 \mathbf{W}_2^T (\mathbf{W}_2 \mathbf{W}_1 \mathbf{X} - \mathbf{Y}) \mathbf{X}^T \mathbf{X} - \mathbf{I}_{N_0} (\mathbf{W}_2 \mathbf{W}_1 \mathbf{X} - \mathbf{Y}) \mathbf{X}^T \mathbf{W}_1^T \mathbf{W}_1 \mathbf{X} \quad (62)$$

$$\Leftrightarrow \frac{d}{dt} f(\mathbf{X}) = - \left(\underbrace{\mathbf{W}_2 \mathbf{W}_2^T \otimes \mathbf{X}^T \mathbf{X} + \mathbf{I} \otimes \mathbf{X}^T \mathbf{W}_1^T \mathbf{W}_1 \mathbf{X}}_{\text{NTK}} \right) \text{vec}_r(\mathbf{W}_2 \mathbf{W}_1 \mathbf{X} - \mathbf{Y}) \quad (63)$$

$$= - \left([\mathbf{W}_2 \otimes \mathbf{X}^T, \mathbf{I} \otimes \mathbf{X}^T \mathbf{W}_1^T] [\mathbf{W}_2 \otimes \mathbf{X}^T, \mathbf{I} \otimes \mathbf{X}^T \mathbf{W}_1^T]^T \right) \text{vec}_r \left(\frac{\partial \mathcal{L}}{\partial F} \right) \quad (64)$$

$$= - \left([\nabla_{\mathbf{W}_1} f, \nabla_{\mathbf{W}_2} f] [\nabla_{\mathbf{W}_1} f, \nabla_{\mathbf{W}_2} f]^T \right) \frac{\partial \mathcal{L}}{\partial f} \quad (65)$$

$$= - (\nabla_{\theta} f \nabla_{\theta} f^T) \frac{\partial \mathcal{L}}{\partial f}, \quad (66)$$

where $[A, B]$ denotes concatenation.

Appendix C. Exact learning dynamics with prior knowledge

C.1. Proof of theorem 3.1

In the following, we prove that equation (11) is in fact a solution to the matrix Riccati equation arising from gradient flow (equation (39)). We prove the theorem by directly substituting our solution for $\mathbf{Q}\mathbf{Q}^T(t)$ into the matrix Riccati equation.

C.1.1. Unequal input-output dimension. We start with the following equation

$$\begin{aligned} \mathbf{Q}\mathbf{Q}^T(t) &= \underbrace{\left[\mathbf{O}e^{\Lambda \frac{t}{\tau}} \mathbf{O}^T + 2\mathbf{M}\mathbf{M}^T \right]}_{\mathbf{L}} \mathbf{Q}(0) \\ &\quad \times \underbrace{\left[\mathbf{I} + \frac{1}{2} \mathbf{Q}(0)^T \left(\mathbf{O} \left(e^{2\Lambda \frac{t}{\tau}} - \mathbf{I} \right) \Lambda^{-1} \mathbf{O}^T + 4 \frac{t}{\tau} \mathbf{M}\mathbf{M}^T \right) \mathbf{Q}(0) \right]^{-1}}_{\mathbf{C}^{-1}} \\ &\quad \times \underbrace{\mathbf{Q}(0)^T \left[\mathbf{O}e^{\Lambda \frac{t}{\tau}} \mathbf{O}^T + 2\mathbf{M}\mathbf{M}^T \right]}_{\mathbf{R}} \\ &= \mathbf{L}\mathbf{C}^{-1}\mathbf{R}, \end{aligned} \tag{67}$$

$$\tag{68}$$

which is identical to equation (11) in the main text, as we verify in section C.2 (by reversing the derivation from equation (152) to equation (128)). Substituting our solution into the matrix Riccati equation then yields

$$\tau \frac{d}{dt} \mathbf{Q}\mathbf{Q}^T = \mathbf{F}\mathbf{Q}\mathbf{Q}^T + \mathbf{Q}\mathbf{Q}^T\mathbf{F} - (\mathbf{Q}\mathbf{Q}^T)^2 \tag{69}$$

$$\Rightarrow \tau \frac{d}{dt} \mathbf{L}\mathbf{C}^{-1}\mathbf{R} \stackrel{?}{=} \mathbf{F}\mathbf{L}\mathbf{C}^{-1}\mathbf{R} + \mathbf{L}\mathbf{C}^{-1}\mathbf{R}\mathbf{F} - \mathbf{L}\mathbf{C}^{-1}\mathbf{R}\mathbf{L}\mathbf{C}^{-1}\mathbf{R}. \tag{70}$$

Next, we note that

$$\mathbf{O}^T \mathbf{O} = \frac{1}{\sqrt{2}} \begin{bmatrix} \tilde{\mathbf{V}} & \tilde{\mathbf{V}} \\ \tilde{\mathbf{U}} & -\tilde{\mathbf{U}} \end{bmatrix}^T \frac{1}{\sqrt{2}} \begin{bmatrix} \tilde{\mathbf{V}} & \tilde{\mathbf{V}} \\ \tilde{\mathbf{U}} & -\tilde{\mathbf{U}} \end{bmatrix} = \mathbf{I}, \tag{71}$$

$$\mathbf{O}^T \mathbf{M} = \frac{1}{\sqrt{2}} \begin{bmatrix} \tilde{\mathbf{V}}^T & \tilde{\mathbf{U}}^T \\ \tilde{\mathbf{V}}^T & -\tilde{\mathbf{U}}^T \end{bmatrix} \frac{1}{\sqrt{2}} \begin{bmatrix} \tilde{\mathbf{V}}_{\perp} \\ \tilde{\mathbf{U}}_{\perp} \end{bmatrix} \tag{72}$$

$$= \frac{1}{2} \begin{bmatrix} \tilde{\mathbf{V}}^T \tilde{\mathbf{V}}_{\perp} + \tilde{\mathbf{U}}^T \tilde{\mathbf{U}}_{\perp} \\ \tilde{\mathbf{V}}^T \tilde{\mathbf{V}}_{\perp} - \tilde{\mathbf{U}}^T \tilde{\mathbf{U}}_{\perp} \end{bmatrix} \tag{73}$$

$$= \mathbf{0} \tag{74}$$

and

$$\mathbf{M}^T \mathbf{O} = \frac{1}{\sqrt{2}} \begin{bmatrix} \tilde{\mathbf{V}}_{\perp}^T & \tilde{\mathbf{U}}_{\perp}^T \end{bmatrix} \frac{1}{\sqrt{2}} \begin{bmatrix} \tilde{\mathbf{V}} & \tilde{\mathbf{V}} \\ \tilde{\mathbf{U}} & -\tilde{\mathbf{U}} \end{bmatrix} \quad (75)$$

$$= \frac{1}{2} \begin{bmatrix} \tilde{\mathbf{V}}_{\perp}^T \tilde{\mathbf{V}} + \tilde{\mathbf{U}}_{\perp}^T \tilde{\mathbf{U}} \\ \tilde{\mathbf{V}}_{\perp}^T \tilde{\mathbf{V}} - \tilde{\mathbf{U}}_{\perp}^T \tilde{\mathbf{U}} \end{bmatrix} \quad (76)$$

$$= \mathbf{0}. \quad (77)$$

Then, using the chain rule $\partial(\mathbf{AB}) = (\partial\mathbf{A})\mathbf{B} + \mathbf{A}(\partial\mathbf{B})$ and the identities

$$\frac{d}{dt}(\mathbf{A}^{-1}) = \mathbf{A}^{-1} \left(\frac{d}{dt} \mathbf{A} \right) \mathbf{A}^{-1} \quad \text{and} \quad \frac{d}{dt}(e^{t\mathbf{A}}) = \mathbf{A}e^{t\mathbf{A}} = e^{t\mathbf{A}}\mathbf{A} \quad (78)$$

we get

$$\tau \frac{d}{dt} \mathbf{Q}\mathbf{Q}^T = \tau \frac{d}{dt} (\mathbf{L}\mathbf{C}^{-1}\mathbf{R}) \quad (79)$$

$$= \tau \left(\frac{d}{dt} \mathbf{L} \right) \mathbf{C}^{-1}\mathbf{R} + \tau \mathbf{L} \left(\frac{d}{dt} \mathbf{C}^{-1}\mathbf{R} \right) \quad (80)$$

$$= \tau \left(\frac{d}{dt} \mathbf{L} \right) \mathbf{C}^{-1}\mathbf{R} + \tau \mathbf{L}\mathbf{C}^{-1} \left(\frac{d}{dt} \mathbf{R} \right) + \tau \mathbf{L} \left(\frac{d}{dt} \mathbf{C}^{-1} \right) \mathbf{R}, \quad (81)$$

with

$$\tau \left(\frac{d}{dt} \mathbf{L} \right) \mathbf{C}^{-1}\mathbf{R} = \tau \mathbf{O} \frac{1}{\tau} \Lambda e^{\Lambda \frac{t}{\tau}} \mathbf{O}^T \mathbf{Q}(0) \mathbf{C}^{-1}\mathbf{R} \quad (82)$$

$$= \mathbf{O} \Lambda e^{\Lambda \frac{t}{\tau}} \mathbf{O}^T \mathbf{Q}(0) \mathbf{C}^{-1}\mathbf{R} \quad (83)$$

$$= \left[\mathbf{O} \Lambda \mathbf{O}^T \mathbf{O} e^{\Lambda \frac{t}{\tau}} \mathbf{O}^T \mathbf{Q}(0) + 2 \mathbf{O} \Lambda \underbrace{\mathbf{O}^T \mathbf{M} \mathbf{M}^T \mathbf{Q}(0)}_{\mathbf{0}} \right] \mathbf{C}^{-1}\mathbf{R} \quad (84)$$

$$= \mathbf{F}\mathbf{L}\mathbf{C}^{-1}\mathbf{R}, \quad (85)$$

$$\tau \mathbf{L}\mathbf{C}^{-1} \left(\frac{d}{dt} \mathbf{R} \right) = \tau \mathbf{L}\mathbf{C}^{-1} \mathbf{Q}(0)^T \mathbf{O} \frac{1}{\tau} e^{\Lambda \frac{t}{\tau}} \Lambda \mathbf{O}^T \quad (86)$$

$$= \mathbf{L}\mathbf{C}^{-1} \mathbf{Q}(0)^T \mathbf{O} e^{\Lambda \frac{t}{\tau}} \Lambda \mathbf{O}^T \quad (87)$$

$$= \mathbf{L}\mathbf{C}^{-1} \left[\mathbf{Q}(0)^T \mathbf{O} e^{\Lambda \frac{t}{\tau}} \mathbf{O}^T \mathbf{O} \Lambda \mathbf{O}^T + 2 \mathbf{Q}(0)^T \mathbf{M} \underbrace{\mathbf{M}^T \mathbf{O} \Lambda \mathbf{O}^T}_{\mathbf{0}} \right] \quad (88)$$

$$= \mathbf{L}\mathbf{C}^{-1}\mathbf{R}\mathbf{F} \quad (89)$$

and

$$\tau \mathbf{L} \left(\frac{d}{dt} \mathbf{C}^{-1} \right) \mathbf{R} = -\tau \mathbf{L} \mathbf{C}^{-1} \left(\frac{d}{dt} \mathbf{C} \right) \mathbf{C}^{-1} \mathbf{R} \quad (90)$$

$$= -\mathbf{L} \mathbf{C}^{-1} \left[\tau \frac{1}{2} \mathbf{Q}(0)^T \mathbf{O} 2 \frac{1}{\tau} e^{2\Lambda \frac{t}{\tau}} \Lambda \Lambda^{-1} \mathbf{O}^T \mathbf{Q}(0) \right. \quad (91)$$

$$\left. + \tau \frac{1}{2} \mathbf{Q}(0)^T 4 \frac{1}{\tau} \mathbf{M} \mathbf{M}^T \mathbf{Q}(0) \right] \mathbf{C}^{-1} \mathbf{R}$$

$$= -\mathbf{L} \mathbf{C}^{-1} \left[\mathbf{Q}(0)^T \mathbf{O} e^{2\Lambda \frac{t}{\tau}} \mathbf{O}^T \mathbf{Q}(0) + 2 \mathbf{Q}(0)^T \mathbf{M} \mathbf{M}^T \mathbf{Q}(0) \right] \mathbf{C}^{-1} \mathbf{R} \quad (92)$$

$$= -\mathbf{L} \mathbf{C}^{-1} \left[\mathbf{Q}(0)^T \mathbf{O} e^{\Lambda \frac{t}{\tau}} \mathbf{O}^T \mathbf{O} e^{\Lambda \frac{t}{\tau}} \mathbf{O}^T \mathbf{Q}(0) \right. \quad (93)$$

$$\left. + 2 \mathbf{Q}(0)^T \mathbf{O} e^{\Lambda \frac{t}{\tau}} \underbrace{\mathbf{O}^T \mathbf{M} \mathbf{M}^T}_{\mathbf{0}} \mathbf{Q}(0) \right.$$

$$\left. + 2 \mathbf{Q}(0)^T \mathbf{M} \mathbf{M}^T \underbrace{\mathbf{O} e^{\Lambda \frac{t}{\tau}} \mathbf{O}^T}_{\mathbf{0}} \mathbf{Q}(0) \right. \quad (94)$$

$$\left. + 4 \mathbf{Q}(0)^T \mathbf{M} \mathbf{M}^T \mathbf{M} \mathbf{M}^T \mathbf{Q}(0) \right] \mathbf{C}^{-1} \mathbf{R}$$

$$= -\mathbf{L} \mathbf{C}^{-1} \mathbf{R} \mathbf{L} \mathbf{C}^{-1} \mathbf{R}.$$

Finally, substituting equations (82), (86) and (90) into the left hand side of equation (70) proves equality. □

C.1.2. Equal input-output dimension. In the case of equal input-output dimensions $\tilde{\mathbf{U}}_{\perp} = \tilde{\mathbf{V}}_{\perp} = 0$ equation (67) reduces to

$$\mathbf{Q} \mathbf{Q}^T(t) = \underbrace{\mathbf{O} e^{\Lambda \frac{t}{\tau}} \mathbf{O}^T}_{\mathbf{L}} \mathbf{Q}(0) \quad (95)$$

$$\times \underbrace{\left[\mathbf{I} + \frac{1}{2} \mathbf{Q}(0)^T \mathbf{O} e^{2\Lambda \frac{t}{\tau}} \Lambda^{-1} \mathbf{O}^T \mathbf{Q}(0) - \frac{1}{2} \mathbf{Q}(0)^T \mathbf{O} \Lambda^{-1} \mathbf{O}^T \mathbf{Q}(0) \right]^{-1}}_{\mathbf{C}^{-1}}$$

$$\times \underbrace{\mathbf{Q}(0)^T \mathbf{O} e^{\Lambda \frac{t}{\tau}} \mathbf{O}^T}_{\mathbf{R}} = \mathbf{L} \mathbf{C}^{-1} \mathbf{R}. \quad (96)$$

Therefore, analogously to the proof for unequal input-output dimensions, it follows that

$$\tau \frac{d}{dt} \mathbf{Q} \mathbf{Q}^T = \tau \frac{d}{dt} \mathbf{L} \mathbf{C}^{-1} \mathbf{R} \quad (97)$$

$$= \tau \left(\frac{d}{dt} \mathbf{L} \right) \mathbf{C}^{-1} \mathbf{R} + \tau \mathbf{L} \left(\frac{d}{dt} \mathbf{C}^{-1} \mathbf{R} \right) \quad (98)$$

$$= \tau \left(\frac{d}{dt} \mathbf{L} \right) \mathbf{C}^{-1} \mathbf{R} + \tau \mathbf{L} \mathbf{C}^{-1} \left(\frac{d}{dt} \mathbf{R} \right) + \tau \mathbf{L} \left(\frac{d}{dt} \mathbf{C}^{-1} \right) \mathbf{R}, \quad (99)$$

with

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

$$= \mathbf{O} \mathbf{\Lambda} \mathbf{O}^T \mathbf{O} e^{\mathbf{\Lambda} \frac{t}{\tau}} \mathbf{O}^T \mathbf{Q}(0) \mathbf{C}^{-1} \mathbf{R} \quad (101)$$

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

$$\tau \mathbf{L} \mathbf{C}^{-1} \left(\frac{d}{dt} \mathbf{R} \right) = \tau \mathbf{L} \mathbf{C}^{-1} \mathbf{Q}(0)^T \mathbf{O} \frac{1}{\tau} e^{\mathbf{\Lambda} \frac{t}{\tau}} \mathbf{\Lambda} \mathbf{O}^T \quad (103)$$

$$= \mathbf{L} \mathbf{C}^{-1} \mathbf{Q}(0)^T \mathbf{O} e^{\mathbf{\Lambda} \frac{t}{\tau}} \mathbf{O}^T \mathbf{O} \mathbf{\Lambda} \mathbf{O}^T \quad (104)$$

$$= \mathbf{L} \mathbf{C}^{-1} \mathbf{R} \mathbf{F}, \quad (105)$$

and

$$\tau \mathbf{L} \left(\frac{d}{dt} \mathbf{C}^{-1} \mathbf{R} \right) = -\tau \mathbf{L} \mathbf{C}^{-1} \left(\frac{d}{dt} \mathbf{C} \right) \mathbf{C}^{-1} \mathbf{R} \quad (106)$$

$$= -\tau \mathbf{L} \mathbf{C}^{-1} \left(\frac{1}{2} \mathbf{Q}(0)^T \mathbf{O} e^{2\mathbf{\Lambda} \frac{t}{\tau}} \frac{2}{\tau} \mathbf{\Lambda} \mathbf{\Lambda}^{-1} \mathbf{O}^T \mathbf{Q}(0) \right) \mathbf{C}^{-1} \mathbf{R} \quad (107)$$

$$= -\tau \mathbf{L} \mathbf{C}^{-1} \mathbf{Q}(0)^T \mathbf{O} e^{\mathbf{\Lambda} \frac{t}{\tau}} \mathbf{O}^T \mathbf{O} e^{\mathbf{\Lambda} \frac{t}{\tau}} \mathbf{Q}(0) \mathbf{C}^{-1} \mathbf{R} \quad (108)$$

$$= -\mathbf{L} \mathbf{C}^{-1} \mathbf{R} \mathbf{L} \mathbf{C}^{-1} \mathbf{R}. \quad (109)$$

Finally, substituting equations (100), (103) and (106) into the left hand side of equation (70) proves equality. \square

C.2. Derivation of the exact learning dynamics

In the following, we outline how the solution to the matrix Riccati equation can be acquired. Let the input and output dimension of a two-layer linear network (equation (1)) be denoted by N_i and N_o respectively. Further, let $N_m = \min(N_i, N_o)$ denote the smaller one of the two. The compact singular value decomposition of the initial network function and the input-output correlation of the task is then

$$\text{SVD}(\mathbf{W}_2(0) \mathbf{W}_1(0)) = \mathbf{U} \mathbf{S} \mathbf{V}^T \text{ and } \text{SVD}(\tilde{\Sigma}^{yx}) = \tilde{\mathbf{U}} \tilde{\mathbf{S}} \tilde{\mathbf{V}}^T. \quad (110)$$

Here, \mathbf{U} and $\tilde{\mathbf{U}} \in \mathbb{R}^{N_o \times N_m}$ denote the left singular vectors, \mathbf{S} and $\tilde{\mathbf{S}} \in \mathbb{R}^{N_m \times N_m}$ the square matrix with ordered, non-zero eigenvalues on its diagonal and \mathbf{V} and $\tilde{\mathbf{V}} \in \mathbb{R}^{N_i \times N_m}$ the

corresponding right singular vectors. Please note that when using compact singular value decomposition, in the case of unequal input-output dimensions ($N_i \neq N_o$) the right and left singular vectors are not generally square and orthonormal.

More specifically, in the case of $N_i < N_o$, $\tilde{\mathbf{U}}^T \tilde{\mathbf{U}} = \tilde{\mathbf{V}}^T \tilde{\mathbf{V}} = \tilde{\mathbf{V}} \tilde{\mathbf{V}}^T = \mathbf{I} \in \mathbb{R}^{N_i \times N_i}$ but $\tilde{\mathbf{U}} \tilde{\mathbf{U}}^T \neq \mathbf{I} \in \mathbb{R}^{N_o \times N_o}$. In this case, we use $\tilde{\mathbf{U}}_{\perp} \in \mathbb{R}^{N_o \times (N_o - N_i)}$ to denote the matrix that contains orthogonal column vectors such that the concatenation $[\tilde{\mathbf{U}} \ \tilde{\mathbf{U}}_{\perp}]$ is orthonormal and $\tilde{\mathbf{V}}_{\perp} \in \mathbb{R}^{N_i \times (N_o - N_i)}$ to denote a matrix of zeros.

Conversely, in the case of $N_i > N_o$, $\tilde{\mathbf{U}} \tilde{\mathbf{U}}^T = \tilde{\mathbf{U}}^T \tilde{\mathbf{U}} = \tilde{\mathbf{V}}^T \tilde{\mathbf{V}} = \mathbf{I} \in \mathbb{R}^{N_o \times N_o}$ but $\tilde{\mathbf{V}}^T \tilde{\mathbf{V}} \neq \mathbf{I} \in \mathbb{R}^{N_i \times N_i}$ and we define $\tilde{\mathbf{V}}_{\perp} \in \mathbb{R}^{N_i \times (N_i - N_o)}$ such that $[\tilde{\mathbf{V}} \ \tilde{\mathbf{V}}_{\perp}]$ is orthonormal and $\tilde{\mathbf{U}}_{\perp} \in \mathbb{R}^{N_o \times (N_o - N_i)}$ to denote a matrix of zeros.

C.2.1. Inverse and matrix exponential of \mathbf{F} . The solution to the matrix Riccati equation as provided by Fukumizu (1998) requires calculation of the inverse \mathbf{F}^{-1} and the matrix exponential $e^{\mathbf{F}\tau}$. To this end, we diagonalise \mathbf{F} by completing its basis by incorporating zero eigenvalues as illustrated below

$$\mathbf{F} = \begin{bmatrix} 0 & \tilde{\mathbf{V}} \tilde{\mathbf{S}} \tilde{\mathbf{U}}^T \\ \tilde{\mathbf{U}} \tilde{\mathbf{S}} \tilde{\mathbf{V}}^T & 0 \end{bmatrix} \tag{111}$$

$$= \frac{1}{\sqrt{2}} \begin{bmatrix} \tilde{\mathbf{V}} & \tilde{\mathbf{V}} & \sqrt{2} \tilde{\mathbf{V}}_{\perp} \\ \tilde{\mathbf{U}} & -\tilde{\mathbf{U}} & \sqrt{2} \tilde{\mathbf{U}}_{\perp} \end{bmatrix} \begin{bmatrix} \tilde{\mathbf{S}} & 0 & 0 \\ 0 & -\tilde{\mathbf{S}} & 0 \\ 0 & 0 & 0 \end{bmatrix} \frac{1}{\sqrt{2}} \begin{bmatrix} \tilde{\mathbf{V}} & \tilde{\mathbf{V}} & \sqrt{2} \tilde{\mathbf{V}}_{\perp} \\ \tilde{\mathbf{U}} & -\tilde{\mathbf{U}} & \sqrt{2} \tilde{\mathbf{U}}_{\perp} \end{bmatrix}^T \tag{112}$$

$$= \mathbf{P} \mathbf{\Gamma} \mathbf{P}^T. \tag{113}$$

Note that $\mathbf{P}^T \mathbf{P} = \mathbf{P} \mathbf{P}^T = \mathbf{I}$ and therefore $\mathbf{P}^T = \mathbf{P}^{-1}$. We then use the diagonalisation of \mathbf{F} to rewrite the matrix exponential

$$e^{\mathbf{F}\tau} = \mathbf{P} e^{\mathbf{\Gamma} \tau} \mathbf{P}^T \tag{114}$$

$$= \frac{1}{\sqrt{2}} \begin{bmatrix} \tilde{\mathbf{V}} & \tilde{\mathbf{V}} & \sqrt{2} \tilde{\mathbf{V}}_{\perp} \\ \tilde{\mathbf{U}} & -\tilde{\mathbf{U}} & \sqrt{2} \tilde{\mathbf{U}}_{\perp} \end{bmatrix} \begin{bmatrix} e^{\tilde{\mathbf{S}}\tau} & 0 & 0 \\ 0 & e^{-\tilde{\mathbf{S}}\tau} & 0 \\ 0 & 0 & e^0 \end{bmatrix} \frac{1}{\sqrt{2}} \begin{bmatrix} \tilde{\mathbf{V}} & \tilde{\mathbf{V}} & \sqrt{2} \tilde{\mathbf{V}}_{\perp} \\ \tilde{\mathbf{U}} & -\tilde{\mathbf{U}} & \sqrt{2} \tilde{\mathbf{U}}_{\perp} \end{bmatrix}^T \tag{115}$$

$$= \frac{1}{2} \begin{bmatrix} \tilde{\mathbf{V}} e^{\tilde{\mathbf{S}}\tau} \tilde{\mathbf{V}}^T + \tilde{\mathbf{V}} e^{-\tilde{\mathbf{S}}\tau} \tilde{\mathbf{V}}^T + 2 \tilde{\mathbf{V}}_{\perp} \tilde{\mathbf{V}}_{\perp}^T & \tilde{\mathbf{V}} e^{\tilde{\mathbf{S}}\tau} \tilde{\mathbf{U}}^T - \tilde{\mathbf{V}} e^{-\tilde{\mathbf{S}}\tau} \tilde{\mathbf{U}}^T + 2 \tilde{\mathbf{V}}_{\perp} \tilde{\mathbf{U}}_{\perp}^T \\ \tilde{\mathbf{U}} e^{\tilde{\mathbf{S}}\tau} \tilde{\mathbf{V}}^T - \tilde{\mathbf{U}} e^{-\tilde{\mathbf{S}}\tau} \tilde{\mathbf{V}}^T + 2 \tilde{\mathbf{U}}_{\perp} \tilde{\mathbf{V}}_{\perp}^T & \tilde{\mathbf{U}} e^{\tilde{\mathbf{S}}\tau} \tilde{\mathbf{U}}^T - \tilde{\mathbf{U}} e^{-\tilde{\mathbf{S}}\tau} \tilde{\mathbf{U}}^T + 2 \tilde{\mathbf{U}}_{\perp} \tilde{\mathbf{U}}_{\perp}^T \end{bmatrix} \tag{116}$$

$$= \frac{1}{\sqrt{2}} \begin{bmatrix} \tilde{\mathbf{V}} & \tilde{\mathbf{V}} \\ \tilde{\mathbf{U}} & -\tilde{\mathbf{U}} \end{bmatrix} \begin{bmatrix} e^{\tilde{\mathbf{S}}\tau} & 0 \\ 0 & e^{-\tilde{\mathbf{S}}\tau} \end{bmatrix} \frac{1}{\sqrt{2}} \begin{bmatrix} \tilde{\mathbf{V}} & \tilde{\mathbf{V}} \\ \tilde{\mathbf{U}} & -\tilde{\mathbf{U}} \end{bmatrix}^T + 2 \frac{1}{\sqrt{2}} \begin{bmatrix} \tilde{\mathbf{V}}_{\perp} \\ \tilde{\mathbf{U}}_{\perp} \end{bmatrix} \frac{1}{\sqrt{2}} \begin{bmatrix} \tilde{\mathbf{V}}_{\perp} \\ \tilde{\mathbf{U}}_{\perp} \end{bmatrix}^T \tag{117}$$

$$= \mathbf{O} e^{\Lambda \tau} \mathbf{O} + 2 \mathbf{M} \mathbf{M}^T. \tag{118}$$

As the inverse $\mathbf{F}^{-1} = \mathbf{P} \mathbf{\Gamma}^{-1} \mathbf{P}^T$ is not well defined for a $\mathbf{\Gamma}$ with zero eigenvalues. We study eigenvalues of value zero by analysing the limiting behaviour of

$$e^{\mathbf{F}\tau} \mathbf{F}^{-1} e^{\mathbf{F}\tau} - \mathbf{F}^{-1} \tag{119}$$

for a single mode

$$\lim_{\epsilon \rightarrow 0} \left[e^{\frac{\epsilon t}{\tau}} \frac{1}{\epsilon} e^{\frac{\epsilon t}{\tau}} - \frac{1}{\epsilon} \right] = \lim_{\epsilon \rightarrow 0} \left[\frac{e^{\frac{2\epsilon t}{\tau}} - 1}{\epsilon} \right] \tag{120}$$

$$\xrightarrow{\text{L'Hospital}} \lim_{\epsilon \rightarrow 0} \left[\frac{\frac{\partial}{\partial \epsilon} \left(e^{\frac{2\epsilon t}{\tau}} - 1 \right)}{\frac{\partial}{\partial \epsilon} \epsilon} \right] \tag{121}$$

$$= \lim_{\epsilon \rightarrow 0} 2 \frac{t}{\tau} e^{\frac{2\epsilon t}{\tau}} \tag{122}$$

$$= 2 \frac{t}{\tau}. \tag{123}$$

which reveals the time dependent contribution of zero eigenvalues. Thus

$$e^{\mathbf{F} \frac{t}{\tau}} \mathbf{F}^{-1} e^{\mathbf{F} \frac{t}{\tau}} - \mathbf{F}^{-1} = \mathbf{O} e^{\Lambda \frac{t}{\tau}} \mathbf{O}^T \mathbf{O} \Lambda^{-1} \mathbf{O}^T \mathbf{O} e^{\Lambda \frac{t}{\tau}} \mathbf{O}^T - \mathbf{O} \Lambda^{-1} \mathbf{O}^T + 4 \frac{t}{\tau} \mathbf{M} \mathbf{M}^T. \tag{124}$$

We continue by substituting the above results into Fukumizu’s equation

$$\mathbf{Q} \mathbf{Q}^T(t) = \left[\mathbf{O} e^{\Lambda \frac{t}{\tau}} \mathbf{O}^T + 2 \mathbf{M} \mathbf{M}^T \right] \mathbf{Q}(0) \tag{125}$$

$$\begin{aligned} & \times \left[\mathbf{I} + \frac{1}{2} \mathbf{Q}(0)^T \left(\mathbf{O} e^{\Lambda \frac{t}{\tau}} \mathbf{O}^T \mathbf{O} \Lambda^{-1} \mathbf{O}^T \mathbf{O} e^{\Lambda \frac{t}{\tau}} \mathbf{O}^T - \mathbf{O} \Lambda^{-1} \mathbf{O}^T + 4 \frac{t}{\tau} \mathbf{M} \mathbf{M}^T \right) \mathbf{Q}(0) \right]^{-1} \\ & \times \mathbf{Q}(0)^T \left[\mathbf{O} e^{\Lambda \frac{t}{\tau}} \mathbf{O}^T + 2 \mathbf{M} \mathbf{M}^T \right] \\ & = \left[\mathbf{O} e^{\Lambda \frac{t}{\tau}} \mathbf{O}^T + 2 \mathbf{M} \mathbf{M}^T \right] \mathbf{Q}(0) \\ & \times \left[\mathbf{I} + \frac{1}{2} \mathbf{Q}(0)^T \left(\mathbf{O} e^{\Lambda \frac{t}{\tau}} \Lambda^{-1} e^{\Lambda \frac{t}{\tau}} \mathbf{O}^T - \mathbf{O} \Lambda^{-1} \mathbf{O}^T + 4 \frac{t}{\tau} \mathbf{M} \mathbf{M}^T \right) \mathbf{Q}(0) \right]^{-1} \end{aligned} \tag{126}$$

$$\begin{aligned} & \mathbf{Q}(0)^T \left[\mathbf{O} e^{\Lambda \frac{t}{\tau}} \mathbf{O}^T + 2 \mathbf{M} \mathbf{M}^T \right] \\ & = \left[\mathbf{O} e^{\Lambda \frac{t}{\tau}} \mathbf{O}^T + 2 \mathbf{M} \mathbf{M}^T \right] \mathbf{Q}(0) \\ & \times \left[\mathbf{I} + \frac{1}{2} \mathbf{Q}(0)^T \left(\mathbf{O} \left(e^{2\Lambda \frac{t}{\tau}} \Lambda^{-1} - \Lambda^{-1} \right) \mathbf{O}^T + 4 \frac{t}{\tau} \mathbf{M} \mathbf{M}^T \right) \mathbf{Q}(0) \right]^{-1} \end{aligned} \tag{127}$$

$$\begin{aligned} & \times \mathbf{Q}(0)^T \left[\mathbf{O} e^{\Lambda \frac{t}{\tau}} \mathbf{O}^T + 2 \mathbf{M} \mathbf{M}^T \right] \\ & = \left[\mathbf{O} e^{\Lambda \frac{t}{\tau}} \mathbf{O}^T + 2 \mathbf{M} \mathbf{M}^T \right] \mathbf{Q}(0) \\ & \times \left[\mathbf{I} + \frac{1}{2} \mathbf{Q}(0)^T \left(\mathbf{O} \left(e^{2\Lambda \frac{t}{\tau}} - \mathbf{I} \right) \Lambda^{-1} \mathbf{O}^T + 4 \frac{t}{\tau} \mathbf{M} \mathbf{M}^T \right) \mathbf{Q}(0) \right]^{-1} \end{aligned} \tag{128}$$

$$\times \mathbf{Q}(0)^T \left[\mathbf{O} e^{\Lambda \frac{t}{\tau}} \mathbf{O}^T + 2 \mathbf{M} \mathbf{M}^T \right].$$

Then, matrix multiplication on the left side of the equation yields

$$\mathbf{O}e^{\Lambda \frac{t}{\tau}} = \frac{1}{\sqrt{2}} \begin{bmatrix} \tilde{\mathbf{V}} & \tilde{\mathbf{V}} \\ \tilde{\mathbf{U}} & -\tilde{\mathbf{U}} \end{bmatrix} \begin{bmatrix} e^{\tilde{\mathbf{S}} \frac{t}{\tau}} & 0 \\ 0 & e^{-\tilde{\mathbf{S}} \frac{t}{\tau}} \end{bmatrix} \tag{129}$$

$$= \frac{1}{\sqrt{2}} \begin{bmatrix} \tilde{\mathbf{V}}e^{\tilde{\mathbf{S}} \frac{t}{\tau}} & \tilde{\mathbf{V}}e^{-\tilde{\mathbf{S}} \frac{t}{\tau}} \\ \tilde{\mathbf{U}}e^{\tilde{\mathbf{S}} \frac{t}{\tau}} & -\tilde{\mathbf{U}}e^{-\tilde{\mathbf{S}} \frac{t}{\tau}} \end{bmatrix} \tag{130}$$

and

$$\mathbf{O}^T \mathbf{Q}(0) = \frac{1}{\sqrt{2}} \begin{bmatrix} \tilde{\mathbf{V}} & \tilde{\mathbf{V}} \\ \tilde{\mathbf{U}} & -\tilde{\mathbf{U}} \end{bmatrix}^T \begin{bmatrix} \mathbf{V}\sqrt{\mathbf{S}}\mathbf{R}^T \\ \mathbf{U}\sqrt{\mathbf{S}}\mathbf{R}^T \end{bmatrix} \tag{131}$$

$$= \frac{1}{\sqrt{2}} \begin{bmatrix} \tilde{\mathbf{V}}^T \mathbf{V}\sqrt{\mathbf{S}}\mathbf{R}^T + \tilde{\mathbf{U}}^T \mathbf{U}\sqrt{\mathbf{S}}\mathbf{R}^T \\ \tilde{\mathbf{V}}^T \mathbf{V}\sqrt{\mathbf{S}}\mathbf{R}^T - \tilde{\mathbf{U}}^T \mathbf{U}\sqrt{\mathbf{S}}\mathbf{R}^T \end{bmatrix} \tag{132}$$

$$= \frac{1}{\sqrt{2}} \begin{bmatrix} \left(\tilde{\mathbf{V}}^T \mathbf{V} + \tilde{\mathbf{U}}^T \mathbf{U} \right) \sqrt{\mathbf{S}}\mathbf{R}^T \\ \left(\tilde{\mathbf{V}}^T \mathbf{V} - \tilde{\mathbf{U}}^T \mathbf{U} \right) \sqrt{\mathbf{S}}\mathbf{R}^T \end{bmatrix}, \tag{133}$$

such that

$$\mathbf{O}e^{\Lambda \frac{t}{\tau}} \mathbf{O}^T \mathbf{Q}(0) = \frac{1}{2} \begin{bmatrix} \tilde{\mathbf{V}}e^{\tilde{\mathbf{S}} \frac{t}{\tau}} & \tilde{\mathbf{V}}e^{-\tilde{\mathbf{S}} \frac{t}{\tau}} \\ \tilde{\mathbf{U}}e^{\tilde{\mathbf{S}} \frac{t}{\tau}} & -\tilde{\mathbf{U}}e^{-\tilde{\mathbf{S}} \frac{t}{\tau}} \end{bmatrix} \begin{bmatrix} \tilde{\mathbf{V}}^T \mathbf{V}\sqrt{\mathbf{S}}\mathbf{R}^T + \tilde{\mathbf{U}}^T \mathbf{U}\sqrt{\mathbf{S}}\mathbf{R}^T \\ \tilde{\mathbf{V}}^T \mathbf{V}\sqrt{\mathbf{S}}\mathbf{R}^T - \tilde{\mathbf{U}}^T \mathbf{U}\sqrt{\mathbf{S}}\mathbf{R}^T \end{bmatrix} \tag{134}$$

$$= \frac{1}{2} \begin{bmatrix} \tilde{\mathbf{V}} \left(e^{\tilde{\mathbf{S}} \frac{t}{\tau}} \left(\tilde{\mathbf{V}}^T \mathbf{V} + \tilde{\mathbf{U}}^T \mathbf{U} \right) + e^{-\tilde{\mathbf{S}} \frac{t}{\tau}} \left(\tilde{\mathbf{V}}^T \mathbf{V} - \tilde{\mathbf{U}}^T \mathbf{U} \right) \right) \sqrt{\mathbf{S}}\mathbf{R}^T \\ \tilde{\mathbf{U}} \left(e^{\tilde{\mathbf{S}} \frac{t}{\tau}} \left(\tilde{\mathbf{V}}^T \mathbf{V} + \tilde{\mathbf{U}}^T \mathbf{U} \right) - e^{-\tilde{\mathbf{S}} \frac{t}{\tau}} \left(\tilde{\mathbf{V}}^T \mathbf{V} - \tilde{\mathbf{U}}^T \mathbf{U} \right) \right) \sqrt{\mathbf{S}}\mathbf{R}^T \end{bmatrix}. \tag{135}$$

We continue by calculating

$$4\mathbf{M}\mathbf{M}^T \mathbf{Q}(0) = 4 \frac{1}{\sqrt{2}} \begin{bmatrix} \tilde{\mathbf{V}}_{\perp} \\ \tilde{\mathbf{U}}_{\perp} \end{bmatrix} \frac{1}{\sqrt{2}} \begin{bmatrix} \tilde{\mathbf{V}}_{\perp} \\ \tilde{\mathbf{U}}_{\perp} \end{bmatrix}^T \begin{bmatrix} \mathbf{V}\sqrt{\mathbf{S}}\mathbf{R}^T \\ \mathbf{U}\sqrt{\mathbf{S}}\mathbf{R}^T \end{bmatrix} \tag{136}$$

$$= 2 \begin{bmatrix} \tilde{\mathbf{V}}_{\perp} \tilde{\mathbf{V}}_{\perp}^T & \tilde{\mathbf{V}}_{\perp} \tilde{\mathbf{U}}_{\perp}^T \\ \tilde{\mathbf{U}}_{\perp} \tilde{\mathbf{V}}_{\perp}^T & \tilde{\mathbf{U}}_{\perp} \tilde{\mathbf{U}}_{\perp}^T \end{bmatrix} \begin{bmatrix} \mathbf{V}\sqrt{\mathbf{S}}\mathbf{R}^T \\ \mathbf{U}\sqrt{\mathbf{S}}\mathbf{R}^T \end{bmatrix} \tag{137}$$

$$= 2 \begin{bmatrix} \tilde{\mathbf{V}}_{\perp} \tilde{\mathbf{V}}_{\perp}^T & 0 \\ 0 & \tilde{\mathbf{U}}_{\perp} \tilde{\mathbf{U}}_{\perp}^T \end{bmatrix} \begin{bmatrix} \mathbf{V}\sqrt{\mathbf{S}}\mathbf{R}^T \\ \mathbf{U}\sqrt{\mathbf{S}}\mathbf{R}^T \end{bmatrix} \tag{138}$$

$$= 2 \begin{bmatrix} \tilde{\mathbf{V}}_{\perp} \tilde{\mathbf{V}}_{\perp}^T \mathbf{V}\sqrt{\mathbf{S}}\mathbf{R}^T \\ \tilde{\mathbf{U}}_{\perp} \tilde{\mathbf{U}}_{\perp}^T \mathbf{U}\sqrt{\mathbf{S}}\mathbf{R}^T \end{bmatrix} \tag{139}$$

and

$$\frac{1}{2}\mathbf{Q}(0)^T 4\frac{t}{\tau}\mathbf{M}\mathbf{M}^T\mathbf{Q}(0) = \frac{t}{\tau} [\mathbf{R}\sqrt{\mathbf{S}}\mathbf{V}^T\mathbf{R}\sqrt{\mathbf{S}}\mathbf{U}^T] \begin{bmatrix} \tilde{\mathbf{V}}_{\perp}\tilde{\mathbf{V}}_{\perp}^T\mathbf{V}\sqrt{\mathbf{S}}\mathbf{R}^T \\ \tilde{\mathbf{U}}_{\perp}\tilde{\mathbf{U}}_{\perp}^T\mathbf{U}\sqrt{\mathbf{S}}\mathbf{R}^T \end{bmatrix} \quad (140)$$

$$= \frac{t}{\tau} \left[\mathbf{R}\sqrt{\mathbf{S}} \left(\mathbf{V}^T\tilde{\mathbf{V}}_{\perp}\tilde{\mathbf{V}}_{\perp}^T\mathbf{V} + \mathbf{U}^T\tilde{\mathbf{U}}_{\perp}\tilde{\mathbf{U}}_{\perp}^T\mathbf{U} \right) \sqrt{\mathbf{S}}\mathbf{R}^T \right] \quad (141)$$

Next, we define $\mathbf{B} = \mathbf{U}^T\tilde{\mathbf{U}} + \mathbf{V}^T\tilde{\mathbf{V}}$ and $\mathbf{C} = \mathbf{U}^T\tilde{\mathbf{U}} - \mathbf{V}^T\tilde{\mathbf{V}}$ and rewrite the inverse as

$$\left[\mathbf{I} + \frac{1}{2}\mathbf{Q}(0)^T \mathbf{O} \left(e^{2\Lambda\frac{t}{\tau}} - \mathbf{I} \right) \Lambda^{-1} \mathbf{O}^T \mathbf{Q}(0) + 2\frac{t}{\tau}\mathbf{Q}(0)^T \mathbf{M}\mathbf{M}^T\mathbf{Q}(0) \right]^{-1} \quad (142)$$

$$= \left[\mathbf{I} + \frac{1}{4}\mathbf{R}\sqrt{\mathbf{S}} \left(\begin{bmatrix} \mathbf{B} & -\mathbf{C} \end{bmatrix} \left(e^{2\Lambda\frac{t}{\tau}} - \mathbf{I} \right) \Lambda^{-1} \begin{bmatrix} \mathbf{B}^T \\ -\mathbf{C}^T \end{bmatrix} + 4\frac{t}{\tau} \left(\mathbf{V}^T\tilde{\mathbf{V}}_{\perp}\tilde{\mathbf{V}}_{\perp}^T\mathbf{V} + \mathbf{U}^T\tilde{\mathbf{U}}_{\perp}\tilde{\mathbf{U}}_{\perp}^T\mathbf{U} \right) \right) \sqrt{\mathbf{S}}\mathbf{R}^T \right]^{-1}. \quad (143)$$

Working from the centre out, we have

$$\begin{bmatrix} \mathbf{B} & -\mathbf{C} \end{bmatrix} \Lambda^{-1} \begin{bmatrix} \mathbf{B}^T \\ -\mathbf{C}^T \end{bmatrix} = \begin{bmatrix} \mathbf{B} & -\mathbf{C} \end{bmatrix} \begin{bmatrix} \tilde{\mathbf{S}}^{-1} & 0 \\ 0 & -\tilde{\mathbf{S}}^{-1} \end{bmatrix} \begin{bmatrix} \mathbf{B}^T \\ -\mathbf{C}^T \end{bmatrix} \quad (144)$$

$$= \begin{bmatrix} \mathbf{B} & -\mathbf{C} \end{bmatrix} \begin{bmatrix} \tilde{\mathbf{S}}^{-1}\mathbf{B}^T \\ \tilde{\mathbf{S}}^{-1}\mathbf{C}^T \end{bmatrix} \quad (145)$$

$$= \mathbf{B}\tilde{\mathbf{S}}^{-1}\mathbf{B}^T - \mathbf{C}\tilde{\mathbf{S}}^{-1}\mathbf{C}^T \quad (146)$$

and

$$\begin{bmatrix} \mathbf{B} & -\mathbf{C} \end{bmatrix} e^{2\Lambda\frac{t}{\tau}} \Lambda^{-1} \begin{bmatrix} \mathbf{B}^T \\ -\mathbf{C}^T \end{bmatrix} = \begin{bmatrix} \mathbf{B} & -\mathbf{C} \end{bmatrix} \begin{bmatrix} e^{2\tilde{\mathbf{S}}\frac{t}{\tau}}\tilde{\mathbf{S}}^{-1} & 0 \\ 0 & -e^{-2\tilde{\mathbf{S}}\frac{t}{\tau}}\tilde{\mathbf{S}}^{-1} \end{bmatrix} \begin{bmatrix} \mathbf{B}^T \\ -\mathbf{C}^T \end{bmatrix} \quad (147)$$

$$= \begin{bmatrix} \mathbf{B} & -\mathbf{C} \end{bmatrix} \begin{bmatrix} e^{2\tilde{\mathbf{S}}\frac{t}{\tau}}\tilde{\mathbf{S}}^{-1}\mathbf{B}^T \\ e^{-2\tilde{\mathbf{S}}\frac{t}{\tau}}\tilde{\mathbf{S}}^{-1}\mathbf{C}^T \end{bmatrix} \quad (148)$$

$$= \mathbf{B}e^{2\tilde{\mathbf{S}}\frac{t}{\tau}}\tilde{\mathbf{S}}^{-1}\mathbf{B}^T - \mathbf{C}e^{-2\tilde{\mathbf{S}}\frac{t}{\tau}}\tilde{\mathbf{S}}^{-1}\mathbf{C}^T. \quad (149)$$

Finally, using $AB^{-1} = (BA^{-1})^{-1}$ (and $A^{-1}B = (B^{-1}A)^{-1}$) to move terms into the inverse, we rewrite

$$\begin{aligned} \mathbf{Q}\mathbf{Q}^T(t) &= \frac{1}{2} \left[\left(\tilde{\mathbf{V}} \left(e^{\tilde{\mathbf{S}}_t} \mathbf{B}^T - e^{-\tilde{\mathbf{S}}_t} \mathbf{C}^T \right) + 2\tilde{\mathbf{V}}_{\perp} \tilde{\mathbf{V}}_{\perp}^T \mathbf{V} \right) \sqrt{\mathbf{S}} \mathbf{R}^T \right] \\ &\quad \times \left[\left(\tilde{\mathbf{U}} \left(e^{\tilde{\mathbf{S}}_t} \mathbf{B}^T + e^{-\tilde{\mathbf{S}}_t} \mathbf{C}^T \right) + 2\tilde{\mathbf{U}}_{\perp} \tilde{\mathbf{U}}_{\perp}^T \mathbf{U} \right) \sqrt{\mathbf{S}} \mathbf{R}^T \right] \\ &\quad \times \left[\mathbf{I} + \mathbf{R} \sqrt{\mathbf{S}} \left(\frac{1}{4} \mathbf{B} \left(e^{2\tilde{\mathbf{S}}_t} - \mathbf{I} \right) \tilde{\mathbf{S}}^{-1} \mathbf{B}^T - \frac{1}{4} \mathbf{C} \left(e^{-2\tilde{\mathbf{S}}_t} - \mathbf{I} \right) \tilde{\mathbf{S}}^{-1} \mathbf{C}^T \right. \right. \\ &\quad \left. \left. + \frac{t}{\tau} \left(\mathbf{V}^T \tilde{\mathbf{V}}_{\perp} \tilde{\mathbf{V}}_{\perp}^T \mathbf{V} + \mathbf{U}^T \tilde{\mathbf{U}}_{\perp} \tilde{\mathbf{U}}_{\perp}^T \mathbf{U} \right) \right) \sqrt{\mathbf{S}} \mathbf{R}^T \right]^{-1} \end{aligned} \tag{150}$$

$$\begin{aligned} &\frac{1}{2} \left[\left(\tilde{\mathbf{V}} \left(e^{\tilde{\mathbf{S}}_t} \mathbf{B}^T - e^{-\tilde{\mathbf{S}}_t} \mathbf{C}^T \right) + 2\tilde{\mathbf{V}}_{\perp} \tilde{\mathbf{V}}_{\perp}^T \mathbf{V} \right) \sqrt{\mathbf{S}} \mathbf{R}^T \right]^T \\ &= \frac{1}{2} \left[\tilde{\mathbf{V}} \left(e^{\tilde{\mathbf{S}}_t} \mathbf{B}^T - e^{-\tilde{\mathbf{S}}_t} \mathbf{C}^T \right) + 2\tilde{\mathbf{V}}_{\perp} \tilde{\mathbf{V}}_{\perp}^T \mathbf{V} \right] \\ &\quad \times \left[\mathbf{S}^{-1} + \frac{1}{4} \mathbf{B} \left(e^{2\tilde{\mathbf{S}}_t} - \mathbf{I} \right) \tilde{\mathbf{S}}^{-1} \mathbf{B}^T - \frac{1}{4} \mathbf{C} \left(e^{-2\tilde{\mathbf{S}}_t} - \mathbf{I} \right) \tilde{\mathbf{S}}^{-1} \mathbf{C}^T \right. \\ &\quad \left. + \frac{t}{\tau} \left(\mathbf{V}^T \tilde{\mathbf{V}}_{\perp} \tilde{\mathbf{V}}_{\perp}^T \mathbf{V} + \mathbf{U}^T \tilde{\mathbf{U}}_{\perp} \tilde{\mathbf{U}}_{\perp}^T \mathbf{U} \right) \right]^{-1} \end{aligned} \tag{151}$$

$$\begin{aligned} &\frac{1}{2} \left[\tilde{\mathbf{V}} \left(e^{\tilde{\mathbf{S}}_t} \mathbf{B}^T - e^{-\tilde{\mathbf{S}}_t} \mathbf{C}^T \right) + 2\tilde{\mathbf{V}}_{\perp} \tilde{\mathbf{V}}_{\perp}^T \mathbf{V} \right]^T \\ &= \left[\tilde{\mathbf{V}} \left(\mathbf{I} - e^{-\tilde{\mathbf{S}}_t} \mathbf{C}^T (\mathbf{B}^T)^{-1} e^{-\tilde{\mathbf{S}}_t} \right) + 2\tilde{\mathbf{V}}_{\perp} \tilde{\mathbf{V}}_{\perp}^T \mathbf{V} (\mathbf{B}^T)^{-1} e^{-\tilde{\mathbf{S}}_t} \right] \\ &\quad \times \left[\tilde{\mathbf{U}} \left(\mathbf{I} + e^{-\tilde{\mathbf{S}}_t} \mathbf{C}^T (\mathbf{B}^T)^{-1} e^{-\tilde{\mathbf{S}}_t} \right) + 2\tilde{\mathbf{U}}_{\perp} \tilde{\mathbf{U}}_{\perp}^T \mathbf{U} (\mathbf{B}^T)^{-1} e^{-\tilde{\mathbf{S}}_t} \right] \\ &\quad \times \left[4e^{-\tilde{\mathbf{S}}_t} \mathbf{B}^{-1} \mathbf{S}^{-1} (\mathbf{B}^T)^{-1} e^{-\tilde{\mathbf{S}}_t} + \left(\mathbf{I} - e^{-2\tilde{\mathbf{S}}_t} \right) \tilde{\mathbf{S}}^{-1} \right. \\ &\quad \left. - e^{-\tilde{\mathbf{S}}_t} \mathbf{B}^{-1} \mathbf{C} \left(e^{-2\tilde{\mathbf{S}}_t} - \mathbf{I} \right) \tilde{\mathbf{S}}^{-1} \mathbf{C}^T (\mathbf{B}^T)^{-1} e^{-\tilde{\mathbf{S}}_t} \right. \\ &\quad \left. + 4\frac{t}{\tau} e^{-\tilde{\mathbf{S}}_t} \mathbf{B}^{-1} \left(\mathbf{V}^T \tilde{\mathbf{V}}_{\perp} \tilde{\mathbf{V}}_{\perp}^T \mathbf{V} + \mathbf{U}^T \tilde{\mathbf{U}}_{\perp} \tilde{\mathbf{U}}_{\perp}^T \mathbf{U} \right) (\mathbf{B}^T)^{-1} e^{-\tilde{\mathbf{S}}_t} \right]^{-1} \\ &\quad \times \left[\tilde{\mathbf{V}} \left(\mathbf{I} - e^{-\tilde{\mathbf{S}}_t} \mathbf{C}^T (\mathbf{B}^T)^{-1} e^{-\tilde{\mathbf{S}}_t} \right) + 2\tilde{\mathbf{V}}_{\perp} \tilde{\mathbf{V}}_{\perp}^T \mathbf{V} \mathbf{B}^{-T} e^{-\tilde{\mathbf{S}}_t} \right]^T \\ &\quad \times \left[\tilde{\mathbf{U}} \left(\mathbf{I} + e^{-\tilde{\mathbf{S}}_t} \mathbf{C}^T (\mathbf{B}^T)^{-1} e^{-\tilde{\mathbf{S}}_t} \right) + 2\tilde{\mathbf{U}}_{\perp} \tilde{\mathbf{U}}_{\perp}^T \mathbf{U} \mathbf{B}^{-T} e^{-\tilde{\mathbf{S}}_t} \right]^T. \end{aligned} \tag{152}$$

C.3. Proof of theorem 3.2: Limiting behaviour

As training time increases, all terms including a matrix exponential with negative exponent in equation (11) vanish to zero, as $\tilde{\mathbf{S}}$ is a diagonal matrix with entries larger zero

$$\lim_{t \rightarrow \infty} e^{-\tilde{\mathbf{S}} \frac{t}{\tau}} = \mathbf{0}. \tag{153}$$

Therefore, in the temporal limit, equation (11) reduces to

$$\lim_{t \rightarrow \infty} \mathbf{Q}\mathbf{Q}^T(t) = \lim_{t \rightarrow \infty} \begin{bmatrix} \mathbf{W}_1^T \mathbf{W}_1(t) & \mathbf{W}_1^T \mathbf{W}_2^T(t) \\ \mathbf{W}_2 \mathbf{W}_1(t) & \mathbf{W}_2^T \mathbf{W}_2(t) \end{bmatrix} \tag{154}$$

$$= \begin{bmatrix} \tilde{\mathbf{V}} \\ \tilde{\mathbf{U}} \end{bmatrix} [\tilde{\mathbf{S}}^{-1}]^{-1} [\tilde{\mathbf{V}}^T \quad \tilde{\mathbf{U}}^T] \tag{155}$$

$$= \begin{bmatrix} \tilde{\mathbf{V}} \tilde{\mathbf{S}} \tilde{\mathbf{V}}^T & \tilde{\mathbf{V}} \tilde{\mathbf{S}} \tilde{\mathbf{U}}^T \\ \tilde{\mathbf{U}} \tilde{\mathbf{S}} \tilde{\mathbf{V}}^T & \tilde{\mathbf{U}} \tilde{\mathbf{S}} \tilde{\mathbf{U}}^T \end{bmatrix}. \tag{156}$$

□

C.4. Dynamics of $\mathbf{Q}(t)$

The solution for the weights $\mathbf{W}_1(t)$ and $\mathbf{W}_2(t)$ can be derived up to a time varying orthogonal transformation as demonstrated by Yan *et al* (1994).

Under the assumptions of whitened inputs 2.2, zero-balanced weights 2.3, full rank 2.4, and equal input-output dimension, the temporal dynamics of $\mathbf{Q}(t)$ is given as

$$\mathbf{Q}(t) = e^{\mathbf{F} \frac{t}{\tau}} \mathbf{Q}(0) \left[\mathbf{I} + \frac{1}{2} \mathbf{Q}(0)^T \left(e^{\mathbf{F} \frac{t}{\tau}} \mathbf{F}^{-1} e^{\mathbf{F} \frac{t}{\tau}} - \mathbf{F}^{-1} \right) \mathbf{Q}(0) \right]^{-\frac{1}{2}} \mathbf{D}(t). \tag{157}$$

where $\mathbf{D}(t)$ is an orthogonal matrix of size $N_h \times N_h$. From this definition, computing $\mathbf{Q}(t)\mathbf{Q}(t)^T$, we recover equation (45).

Equation (157) shows that the individual weight matrices are not directly described by parts of the $\mathbf{Q}(t)\mathbf{Q}(t)^T$ solution. Instead, they are fixed only up to a time-dependent orthogonal transformation. To verify this, we numerically compute $\mathbf{D}(t)$ as $\mathbf{D}(t) = \mathbf{q}(t)^+ \mathbf{Q}_{\text{sim}}(t)$ where $\mathbf{Q}_{\text{sim}}(t)$ denotes weights obtained from numerical simulations of gradient descent, $+$ denotes the pseudoinverse ($\mathbf{q}^+(t) = (\mathbf{q}^T(t)\mathbf{q}(t))^{-1}\mathbf{q}(t)^T$ where $\mathbf{q}(t)$ is rectangular) and

$$\mathbf{q}(t) = e^{\mathbf{F} \frac{t}{\tau}} \mathbf{Q}(0) \left[\mathbf{I} + \frac{1}{2} \mathbf{Q}(0)^T \left(e^{\mathbf{F} \frac{t}{\tau}} \mathbf{F}^{-1} e^{\mathbf{F} \frac{t}{\tau}} - \mathbf{F}^{-1} \right) \mathbf{Q}(0) \right]^{-\frac{1}{2}}. \tag{158}$$

We numerically show in figure 7(D) right panel that $\mathbf{D}(t)$ generally changes over time. Letting $\mathbf{Q}_d(t)$ denote the estimated $\mathbf{Q}(t)$ using the numerically recovered $\mathbf{D}(t)$, figure 7(D) left and centre panels show that both the dynamics of $\mathbf{Q}_d(t)$ and $\mathbf{Q}_d(t)\mathbf{Q}_d(t)^T$ match the temporal dynamics of the simulation. The small derivation between the

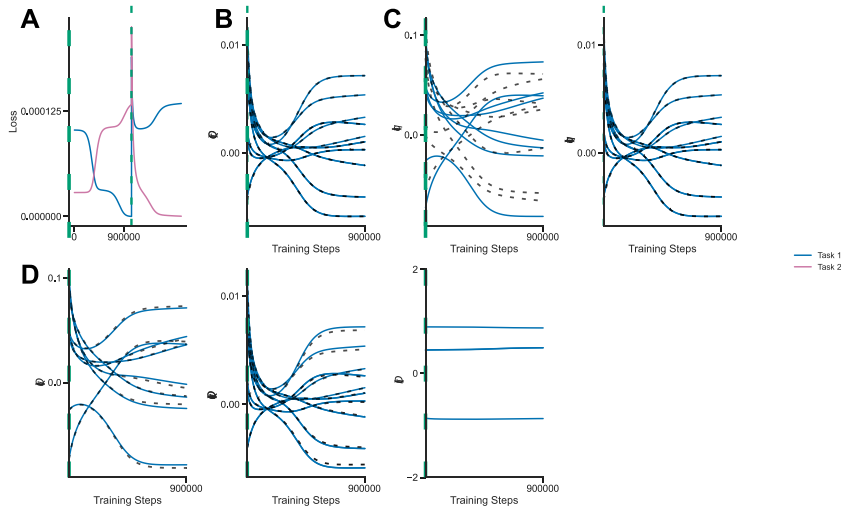


Figure 7. (A) Loss under gradient descent learning two random input-output correlation task with learning rate $\eta = 0,001$ up to precision 1×10^{-7} . The green dotted line marks the time at which the target is switched from task 1 to task 2. (B) Numerical (coloured line) and analytical (black dotted line) temporal dynamics of $\mathbf{Q}\mathbf{Q}^T(t)$ as given by equation (159). (C) Numerical (coloured line) and analytical (black dotted line) temporal dynamics of $\mathbf{q}(t)$ and $\mathbf{q}(t)\mathbf{q}(t)^T$ (158) (D) Temporal dynamics of $\mathbf{D}(t)$. Numerical (coloured line) and analytical (black dotted line) temporal dynamics of $\mathbf{Q}_d(t)\mathbf{Q}_d(t)^T$ and $\mathbf{Q}_d(t)$ as given by equation (157) where (D) was computed numerically.

simulation and the analytical solution for later time points, is due to the imprecision of the pseudoinverse.

In figure 7(C), we report the implementation of equation (158). As expected, the analytical solution does not match the numerical temporal dynamics. However, the solution for $\mathbf{q}(t)\mathbf{q}(t)^T$ recovers the correct dynamics.

Appendix D. Rich and lazy learning regimes and generalisation

Under the assumptions of theorem 3.1, the network function acquires a rich task-specific internal representation at convergence, that is $\mathbf{W}_1^T\mathbf{W}_1 = \tilde{\mathbf{V}}\tilde{\mathbf{S}}\tilde{\mathbf{V}}^T$ and $\mathbf{W}_2\mathbf{W}_2^T = \tilde{\mathbf{U}}\tilde{\mathbf{S}}\tilde{\mathbf{U}}^T$. Therefore, there exist initial states with large zero-balanced weights that lead to rich solutions.

We more quantitatively capture this phenomena in figure 8. We define the error on the internal representation as figure 3 $\|\mathbf{W}_1^T\mathbf{W}_1 - \tilde{\mathbf{V}}\tilde{\mathbf{S}}\tilde{\mathbf{V}}^T\|_F^2$ and $\|\mathbf{W}_2\mathbf{W}_2^T - \tilde{\mathbf{U}}\tilde{\mathbf{S}}\tilde{\mathbf{U}}^T\|_F^2$ for \mathbf{W}_1 and \mathbf{W}_2 respectively. Effectively, we measure the richness of the representation and in turn it is generalisation ability. In figure 8, the error remains zero for increasing gain for any network initialised with zero-balanced weights. In other words, the representation at convergences is rich. In contrast, for random initialisation the error increase consequently with increasing gain. As the network is moving away from the small random weight initialisation, the network converges to lazier representation.

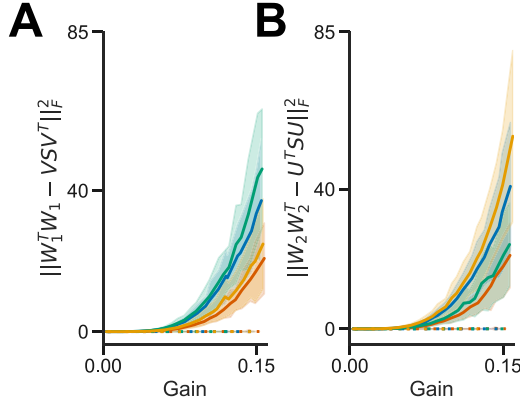


Figure 8. (A) and (B) Mean and standard deviation on the error on the internal representation error defined as in section D for the learning the living kingdom task (figure 6(A)), a random 7×7 matrix (blue), a random 5×7 matrix (yellow), a 7×5 matrix (green), a 8×8 matrix (red). All the task ran were ran with learning rate $\eta = 0.001$ enforcing initial zero-balanced weights 2.3 (dotted line) and breaking the assumption of zero-balanced initial weights 2.3 (line). $N_h = 10$ for all networks.

Appendix E. Decoupling dynamics

E.1. Proof for theorem 5.1

Let the input and output dimension of a two-layer linear network (equation (1)) be equal, i.e. $N_i = N_o$, then equation (11) simplifies to

$$\begin{aligned}
 \mathbf{Q}\mathbf{Q}^T(t) &= \begin{bmatrix} \tilde{\mathbf{V}} \left(\mathbf{I} - e^{-\tilde{\mathbf{S}} \frac{t}{\tau}} \mathbf{C}^T (\mathbf{B}^T)^{-1} e^{-\tilde{\mathbf{S}} \frac{t}{\tau}} \right) \\ \tilde{\mathbf{U}} \left(\mathbf{I} + e^{-\tilde{\mathbf{S}} \frac{t}{\tau}} \mathbf{C}^T (\mathbf{B}^T)^{-1} e^{-\tilde{\mathbf{S}} \frac{t}{\tau}} \right) \end{bmatrix} \\
 &\quad \times \left[4e^{-\tilde{\mathbf{S}} \frac{t}{\tau}} \mathbf{B}^{-1} \mathbf{S}^{-1} (\mathbf{B}^T)^{-1} e^{-\tilde{\mathbf{S}} \frac{t}{\tau}} + (\mathbf{I} - e^{-2\tilde{\mathbf{S}} \frac{t}{\tau}}) \tilde{\mathbf{S}}^{-1} \right. \\
 &\quad \left. - e^{-\tilde{\mathbf{S}} \frac{t}{\tau}} \mathbf{B}^{-1} \mathbf{C} \left(e^{-2\tilde{\mathbf{S}} \frac{t}{\tau}} - \mathbf{I} \right) \tilde{\mathbf{S}}^{-1} \mathbf{C}^T (\mathbf{B}^T)^{-1} e^{-\tilde{\mathbf{S}} \frac{t}{\tau}} \right]^{-1} \\
 &\quad \times \begin{bmatrix} \tilde{\mathbf{V}} \left(\mathbf{I} - e^{-\tilde{\mathbf{S}} \frac{t}{\tau}} \mathbf{C}^T (\mathbf{B}^T)^{-1} e^{-\tilde{\mathbf{S}} \frac{t}{\tau}} \right) \\ \tilde{\mathbf{U}} \left(\mathbf{I} + e^{-\tilde{\mathbf{S}} \frac{t}{\tau}} \mathbf{C}^T (\mathbf{B}^T)^{-1} e^{-\tilde{\mathbf{S}} \frac{t}{\tau}} \right) \end{bmatrix}^T. \tag{159}
 \end{aligned}$$

Further, let the singular value decomposition of the input-output correlation of the task be

$$\text{SVD} \left(\tilde{\Sigma}^{yx} \right) = \tilde{\mathbf{U}} \tilde{\mathbf{S}} \tilde{\mathbf{V}}^T \tag{160}$$

and suppose that the initial state of the network can be written in the form

$$\text{SVD} \left(\mathbf{W}_2(0) \mathbf{W}_1(0) \right) = \mathbf{U} \mathbf{S} \mathbf{V}^T = \tilde{\mathbf{U}} \mathbf{A}(0)^T \mathbf{A}(0) \tilde{\mathbf{V}}^T. \tag{161}$$

First, we note that the initial weights in this setting are not independent of the structure of the target task. In particular,

$$\mathbf{U}\sqrt{\mathbf{S}} = \tilde{\mathbf{U}}\mathbf{A}(0)^{\text{T}} \tag{162}$$

$$\Leftrightarrow \tilde{\mathbf{U}}^{\text{T}}\mathbf{U}\sqrt{\mathbf{S}} = \mathbf{A}(0)^{\text{T}} \tag{163}$$

$$\Leftrightarrow \sqrt{\mathbf{S}}\mathbf{U}^{\text{T}}\tilde{\mathbf{U}} = \mathbf{A}(0) \tag{164}$$

and

$$\sqrt{\mathbf{S}}\mathbf{V}^{\text{T}} = \mathbf{A}(0)\tilde{\mathbf{V}}^{\text{T}} \tag{166}$$

$$\Leftrightarrow \sqrt{\mathbf{S}}\mathbf{V}^{\text{T}}\tilde{\mathbf{V}} = \mathbf{A}(0) \tag{167}$$

and therefore

$$\sqrt{\mathbf{S}}\mathbf{U}^{\text{T}}\tilde{\mathbf{U}} = \sqrt{\mathbf{S}}\mathbf{V}^{\text{T}}\tilde{\mathbf{V}} \tag{168}$$

$$\Leftrightarrow \mathbf{U}\mathbf{V}^{\text{T}} = \tilde{\mathbf{U}}\tilde{\mathbf{V}}^{\text{T}}. \tag{169}$$

This further simplifies the equation, as

$$\mathbf{U}\sqrt{\mathbf{S}} = \tilde{\mathbf{U}}\mathbf{A}(0)^{\text{T}} \tag{170}$$

$$\Leftrightarrow \mathbf{U} = \tilde{\mathbf{U}}\mathbf{A}(0)^{\text{T}}\sqrt{\mathbf{S}}^{-1} \tag{171}$$

and

$$\sqrt{\mathbf{S}}\mathbf{V}^{\text{T}} = \mathbf{A}(0)\tilde{\mathbf{V}}^{\text{T}} \tag{172}$$

$$\Leftrightarrow \mathbf{V}^{\text{T}} = \sqrt{\mathbf{S}}^{-1}\mathbf{A}(0)\tilde{\mathbf{V}}^{\text{T}} \tag{173}$$

$$\Leftrightarrow \mathbf{V} = \tilde{\mathbf{V}}\mathbf{A}(0)^{\text{T}}\sqrt{\mathbf{S}}^{-1}, \tag{174}$$

then recollecting the definition of \mathbf{B} and \mathbf{C} we get

$$\mathbf{B}^{\text{T}} = \tilde{\mathbf{U}}^{\text{T}}\mathbf{U} + \tilde{\mathbf{V}}^{\text{T}}\mathbf{V} \tag{175}$$

$$= \tilde{\mathbf{U}}^{\text{T}}\tilde{\mathbf{U}}\mathbf{A}(0)^{\text{T}}\sqrt{\mathbf{S}}^{-1} + \tilde{\mathbf{V}}^{\text{T}}\tilde{\mathbf{V}}\mathbf{A}(0)^{\text{T}}\sqrt{\mathbf{S}}^{-1} \tag{176}$$

$$= \left(\tilde{\mathbf{U}}^{\text{T}}\tilde{\mathbf{U}} + \tilde{\mathbf{V}}^{\text{T}}\tilde{\mathbf{V}}\right)\mathbf{A}(0)^{\text{T}}\sqrt{\mathbf{S}}^{-1} \tag{177}$$

$$= 2\mathbf{A}(0)^{\text{T}}\sqrt{\mathbf{S}}^{-1} \tag{178}$$

and

$$\mathbf{C}^{\text{T}} = \tilde{\mathbf{U}}^{\text{T}}\mathbf{U} - \tilde{\mathbf{V}}^{\text{T}}\mathbf{V} \tag{179}$$

$$= \left(\tilde{\mathbf{U}}^{\text{T}}\tilde{\mathbf{U}} - \tilde{\mathbf{V}}^{\text{T}}\tilde{\mathbf{V}}\right)\mathbf{A}(0)^{\text{T}}\sqrt{\mathbf{S}}^{-1} \tag{180}$$

$$= 0. \tag{181}$$

Substituting the new values of \mathbf{B} and \mathbf{C} into equation (159) then yields

$$\mathbf{Q}\mathbf{Q}^T(t) = \begin{bmatrix} \tilde{\mathbf{V}} \\ \tilde{\mathbf{U}} \end{bmatrix} \left[4e^{-\tilde{\mathbf{S}}\frac{t}{\tau}} \frac{1}{4} \mathbf{A}(0)^{-1} \sqrt{\mathbf{S}} \mathbf{S}^{-1} \sqrt{\mathbf{S}} \mathbf{A}(0)^{-T} e^{-\tilde{\mathbf{S}}\frac{t}{\tau}} + \left(\mathbf{I} - e^{-2\tilde{\mathbf{S}}\frac{t}{\tau}} \right) \tilde{\mathbf{S}}^{-1} \right]^{-1} \begin{bmatrix} \tilde{\mathbf{V}} \\ \tilde{\mathbf{U}} \end{bmatrix}^T \quad (182)$$

$$= \begin{bmatrix} \tilde{\mathbf{V}} \\ \tilde{\mathbf{U}} \end{bmatrix} \left[e^{-\tilde{\mathbf{S}}\frac{t}{\tau}} \left(\mathbf{A}(0)^T \mathbf{A}(0) \right)^{-1} e^{-\tilde{\mathbf{S}}\frac{t}{\tau}} + \left(\mathbf{I} - e^{-2\tilde{\mathbf{S}}\frac{t}{\tau}} \right) \tilde{\mathbf{S}}^{-1} \right]^{-1} \begin{bmatrix} \tilde{\mathbf{V}} \\ \tilde{\mathbf{U}} \end{bmatrix}^T. \quad (183)$$

Finally, we note that the dynamics can thus be written as

$$\mathbf{Q}\mathbf{Q}^T(t) = \begin{bmatrix} \tilde{\mathbf{V}} \mathbf{A}^T \mathbf{A}(t) \tilde{\mathbf{V}}^T & \tilde{\mathbf{V}} \mathbf{A}^T \mathbf{A}(t) \tilde{\mathbf{U}}^T \\ \tilde{\mathbf{U}} \mathbf{A}^T \mathbf{A}(t) \tilde{\mathbf{V}}^T & \tilde{\mathbf{U}} \mathbf{A}^T \mathbf{A}(t) \tilde{\mathbf{U}}^T \end{bmatrix} \quad (184)$$

where

$$\mathbf{A}^T \mathbf{A}(t) = \left[e^{-\tilde{\mathbf{S}}\frac{t}{\tau}} \left(\mathbf{A}(0)^T \mathbf{A}(0) \right)^{-1} e^{-\tilde{\mathbf{S}}\frac{t}{\tau}} + \left(\mathbf{I} - e^{-2\tilde{\mathbf{S}}\frac{t}{\tau}} \right) \tilde{\mathbf{S}}^{-1} \right]^{-1}. \quad (185)$$

□

E.2. Solution for 2×2 dynamics

We consider small networks with input and output dimension $N_i = 2$ and $N_o = 2$. In this setting, the structure of the weight initialisation and task are encoded in the matrices

$$\mathbf{A}(0)^T \mathbf{A}(0) = \begin{bmatrix} a_1(0) & b(0) \\ b(0) & a_2(0) \end{bmatrix} \text{ and } \tilde{\mathbf{S}} = \begin{bmatrix} s_1 & 0 \\ 0 & s_2 \end{bmatrix}, \quad (186)$$

where the parameters $a_1(0)$ and $a_2(0)$ represent coupling within a singular mode, and $b(0)$ represents counterproductive cross-coupling between different singular modes.

From equation (13), we have

$$\begin{aligned} \mathbf{A}^T \mathbf{A}(t) &= \left[\begin{bmatrix} e^{-\frac{s_1 t}{\tau}} & 0 \\ 0 & e^{-\frac{s_2 t}{\tau}} \end{bmatrix} \begin{bmatrix} a_1(0) & b(0) \\ b(0) & a_2(0) \end{bmatrix}^{-1} \begin{bmatrix} e^{-\frac{s_1 t}{\tau}} & 0 \\ 0 & e^{-\frac{s_2 t}{\tau}} \end{bmatrix} \right. \\ &\quad \left. + \left[\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} - \begin{bmatrix} e^{-\frac{2s_1 t}{\tau}} & 0 \\ 0 & e^{-\frac{2s_2 t}{\tau}} \end{bmatrix} \right] \begin{bmatrix} s_1 & 0 \\ 0 & s_2 \end{bmatrix}^{-1} \right]^{-1} \end{aligned} \quad (187)$$

$$\begin{aligned} &= \left[\frac{1}{a_1(0)a_2(0) - b(0)^2} \begin{bmatrix} e^{-\frac{s_1 t}{\tau}} & 0 \\ 0 & e^{-\frac{s_2 t}{\tau}} \end{bmatrix} \begin{bmatrix} a_2(0) & -b(0) \\ -b(0) & a_1(0) \end{bmatrix} \begin{bmatrix} e^{-\frac{s_1 t}{\tau}} & 0 \\ 0 & e^{-\frac{s_2 t}{\tau}} \end{bmatrix} \right. \\ &\quad \left. + \left[\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} - \begin{bmatrix} e^{-\frac{2s_1 t}{\tau}} & 0 \\ 0 & e^{-\frac{2s_2 t}{\tau}} \end{bmatrix} \right] \begin{bmatrix} \frac{1}{s_1} & 0 \\ 0 & \frac{1}{s_2} \end{bmatrix} \right]^{-1}, \end{aligned} \quad (188)$$

where we use

$$\begin{bmatrix} a & b \\ c & d \end{bmatrix}^{-1} = \frac{1}{ad - bc} \begin{bmatrix} d & -b \\ -c & a \end{bmatrix}. \tag{189}$$

We continue with

$$\begin{aligned} A^T A(t) &= \left[\frac{1}{a_1(0)a_2(0)-b(0)^2} \begin{bmatrix} e^{-\frac{2s_1 t}{\tau}} & 0 \\ 0 & e^{-\frac{2s_2 t}{\tau}} \end{bmatrix} \begin{bmatrix} a_2(0) & -b(0) \\ -b(0) & a_1(0) \end{bmatrix} \begin{bmatrix} e^{-\frac{2s_1 t}{\tau}} & 0 \\ 0 & e^{-\frac{2s_2 t}{\tau}} \end{bmatrix} \right. \\ &\quad \left. + \begin{bmatrix} \frac{1}{s_1} & 0 \\ 0 & \frac{1}{s_2} \end{bmatrix} - \begin{bmatrix} \frac{1}{s_1} e^{-\frac{2s_1 t}{\tau}} & 0 \\ 0 & \frac{1}{s_2} e^{-\frac{2s_2 t}{\tau}} \end{bmatrix} \right]^{-1} \end{aligned} \tag{190}$$

$$\begin{aligned} &= \left[\frac{1}{a_1(0)a_2(0)-b(0)^2} \begin{bmatrix} e^{-\frac{2s_1 t}{\tau}} a_2(0) & -e^{-\frac{s_1 t}{\tau}} b(0) e^{-\frac{s_2 t}{\tau}} \\ -e^{-\frac{s_2 t}{\tau}} b(0) e^{-\frac{s_1 t}{\tau}} & e^{-\frac{2s_2 t}{\tau}} a_1(0) \end{bmatrix} \right. \\ &\quad \left. + \begin{bmatrix} \frac{1}{s_1} & 0 \\ 0 & \frac{1}{s_2} \end{bmatrix} - \begin{bmatrix} \frac{1}{s_1} e^{-\frac{2s_1 t}{\tau}} & 0 \\ 0 & \frac{1}{s_2} e^{-\frac{2s_2 t}{\tau}} \end{bmatrix} \right]^{-1} \end{aligned} \tag{191}$$

$$= \begin{bmatrix} \frac{e^{-\frac{2s_1 t}{\tau}} a_2(0)}{a_1(0)a_2(0)-b(0)^2} + \frac{1}{s_1} - \frac{1}{s_1} e^{-\frac{2s_1 t}{\tau}} & -\frac{e^{-\frac{s_1 t}{\tau}} b(0) e^{-\frac{s_2 t}{\tau}}}{a_1(0)a_2(0)-b(0)^2} \\ -\frac{e^{-\frac{s_2 t}{\tau}} b(0) e^{-\frac{s_1 t}{\tau}}}{a_1(0)a_2(0)-b(0)^2} & \frac{e^{-\frac{2s_2 t}{\tau}} a_1(0)}{a_1(0)a_2(0)-b(0)^2} + \frac{1}{s_2} - \frac{1}{s_2} e^{-\frac{2s_2 t}{\tau}} \end{bmatrix}^{-1}. \tag{192}$$

We use equation (189) and simplify the denominator

$$\begin{aligned} \mathbf{A}^T \mathbf{A}(t) &= \frac{1}{\left(\frac{e^{-\frac{2s_2 t}{\tau}} a_1(0)}{a_1(0)a_2(0)-b(0)^2} + \frac{1}{s_2} - \frac{1}{s_2} e^{-\frac{2s_2 t}{\tau}} \right) \left(\frac{e^{-\frac{2s_1 t}{\tau}} a_2(0)}{a_1(0)a_2(0)-b(0)^2} + \frac{1}{s_1} - \frac{1}{s_1} e^{-\frac{2s_1 t}{\tau}} \right) - \left(-\frac{e^{-\frac{s_2 t}{\tau}} b(0) e^{-\frac{s_1 t}{\tau}}}{a_1(0)a_2(0)-b(0)^2} \right)^2} \\ &\quad \times \begin{bmatrix} \frac{e^{-\frac{2s_2 t}{\tau}} a_1(0)}{a_1(0)a_2(0)-b(0)^2} + \frac{1}{s_2} - \frac{1}{s_2} e^{-\frac{2s_2 t}{\tau}} & \frac{e^{-\frac{s_1 t}{\tau}} b(0) e^{-\frac{s_2 t}{\tau}}}{a_1(0)a_2(0)-b(0)^2} \\ \frac{e^{-\frac{s_2 t}{\tau}} b(0) e^{-\frac{s_1 t}{\tau}}}{a_1(0)a_2(0)-b(0)^2} & \frac{e^{-\frac{2s_1 t}{\tau}} a_2(0)}{a_1(0)a_2(0)-b(0)^2} + \frac{1}{s_1} - \frac{1}{s_1} e^{-\frac{2s_1 t}{\tau}} \end{bmatrix}. \end{aligned} \tag{193}$$

The diagonal element $a_1(t)$ is given as

$$\begin{aligned}
 a_1(t) &= \frac{\frac{e^{-2s_2t}}{a_1(0)a_2(0)-b(0)^2} a_1(0) + \frac{1}{s_2} - \frac{1}{s_2} e^{-\frac{2s_2t}{\tau}}}{\left(\frac{e^{-\frac{2s_2t}{\tau}} a_1(0)}{a_1(0)a_2(0)-b(0)^2} + \frac{1}{s_2} - \frac{1}{s_2} e^{-\frac{2s_2t}{\tau}} \right) \left(\frac{e^{-\frac{2s_1t}{\tau}} a_2(0)}{a_1(0)a_2(0)-b(0)^2} + \frac{1}{s_1} - \frac{1}{s_1} e^{-\frac{2s_1t}{\tau}} \right) - \left(-\frac{e^{-\frac{s_2t}{\tau}} b(0) e^{-\frac{s_1t}{\tau}}}{a_1(0)a_2(0)-b(0)^2} \right)^2}, \tag{194}
 \end{aligned}$$

and interchanging subscripts 1 and 2 yields $a_2(t)$. As a check on this result, by setting $b(0) = 0$ we recover the expression

$$a_1(t) = \frac{a_1(0)}{e^{-\frac{2s_1t}{\tau}} + \frac{a_1(0)}{s_1} \left(1 - e^{-\frac{2s_1t}{\tau}} \right)}, \tag{195}$$

from Saxe *et al* (2019).

We further simplify the denominator to

$$\begin{aligned}
 \mathbf{A}^T \mathbf{A}(t) &= \frac{1}{\frac{1}{a_1(0)a_2(0)-b(0)^2} \left(e^{-\frac{2(s_1+s_2)t}{\tau}} \left(1 - \frac{a_1(0)}{s_1} - \frac{a_2(0)}{s_2} \right) + e^{-\frac{2s_2t}{\tau}} \frac{a_1(0)}{s_1} + e^{-\frac{2s_1t}{\tau}} \frac{a_2(0)}{s_2} \right) + \frac{1}{s_2s_1}} \\
 &\times \begin{bmatrix} \frac{e^{-\frac{2s_2t}{\tau}} a_1(0)}{a_1(0)a_2(0)-b(0)^2} + \frac{1}{s_2} - \frac{1}{s_2} e^{-\frac{2s_2t}{\tau}} & \frac{e^{-\frac{s_1t}{\tau}} b(0) e^{-\frac{s_2t}{\tau}}}{a_1(0)a_2(0)-b(0)^2} \\ \frac{e^{-\frac{s_2t}{\tau}} b(0) e^{-\frac{s_1t}{\tau}}}{a_1(0)a_2(0)-b(0)^2} & \frac{e^{-\frac{2s_1t}{\tau}} a_2(0)}{a_1(0)a_2(0)-b(0)^2} + \frac{1}{s_1} - \frac{1}{s_1} e^{-\frac{2s_1t}{\tau}} \end{bmatrix} \tag{196}
 \end{aligned}$$

E.3. Off-Diagonal decoupling dynamics

We track the decoupling by considering the dynamics of the off-diagonal element $b(t)$.

$$b(t) = \frac{\frac{e^{-\frac{s_2t}{\tau}} b(0) e^{-\frac{s_1t}{\tau}}}{a_1(0)a_2(0)-b(0)^2}}{\frac{1}{a_1(0)a_2(0)-b(0)^2} \left(e^{-\frac{2(s_1+s_2)t}{\tau}} \left(1 - \frac{a_1(0)}{s_1} - \frac{a_2(0)}{s_2} \right) + e^{-\frac{2s_2t}{\tau}} \frac{a_1(0)}{s_1} + e^{-\frac{2s_1t}{\tau}} \frac{a_2(0)}{s_2} \right) + \frac{1}{s_2s_1}}. \tag{197}$$

As t tends to infinity $\lim_{t \rightarrow \infty} b(t) = 0$ the off-diagonal element shrinks to zero.

We can further simplify the off-diagonal to

$$b(t) = \frac{b(0)}{e^{-\frac{(s_1+s_2)t}{\tau}} \left(1 - \frac{a_1(0)}{s_1} - \frac{a_2(0)}{s_2} \right) + e^{\frac{(s_1-s_2)t}{\tau}} \frac{a_1(0)}{s_1} + e^{\frac{(s_2-s_1)t}{\tau}} \frac{a_2(0)}{s_2} + \frac{a_1(0)a_2(0)-b(0)^2}{s_2s_1}}. \tag{198}$$

Equation (198) can exhibit non-monotonic trajectories with transient peaks as shown in figure 4. The qualitative observations for the 2×2 network hold for larger target matrices as shown in figure 9. For large initialisation, the dynamics are exponential.

At intermediate and small initialisation, the maximum of the off-diagonal is reached before the singular mode is fully learned. In the small initialisation scheme, the peak is of negligible size. The respective target matrix for panel (A)–(D), (B)–(E) and (C)–(F) in figure 9 are

$$\text{dense } \begin{bmatrix} 5 & 6 & 3 & 0 & 1 \\ 4, & 1 & 0 & 1 & 2 \\ 3 & 0 & 2 & 4 & 0 \\ 3 & 4 & 0 & 3 & 2 \\ 2 & 0 & 1 & 3 & 4 \end{bmatrix}, \text{ diagonal } \begin{bmatrix} 5 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 2 & 0 & 0 \\ 0 & 0 & 0 & 3 & 0 \\ 0 & 0 & 0 & 0 & 4 \end{bmatrix} \text{ and equal diagonal } \begin{bmatrix} 5 & 0 & 0 & 0 & 0 \\ 0 & 5 & 0 & 0 & 0 \\ 0 & 0 & 5 & 0 & 0 \\ 0 & 0 & 0 & 5 & 0 \\ 0 & 0 & 0 & 0 & 5 \end{bmatrix}.$$

We characterise these dynamics considering the case where $s_1 = s_2 = s$ for the two-by-two solution (i.e. equal diagonal target \mathbf{y}) for which we can compute the time of the peak. In this particular case, we can further simplify the off-diagonal to

$$b(t) = \frac{b(0)}{e^{-\frac{2(s)t}{\tau}} \left(1 - \frac{a_1(0)+a_2(0)}{s}\right) + \frac{a_1(0)+a_2(0)}{s} + \frac{a_1(0)a_2(0)-b(0)^2}{s^2}}. \tag{199}$$

We find the time of the maximum of the off-diagonal elements to be $t_{\text{peak}} = \frac{\tau}{4s} \ln \frac{s(s-a_1(0)-a_2(0))}{a_1(0)a_2(0)-b(0)^2}$.

The presence of a peak in the off-diagonal values, indicates the decoupling, but as shown in figures 4(D)–(F), the peak size is negligible in comparison to the size of the on-diagonal values for small initial weights. This difference is reminiscent of the silent alignment effect described by Atanasov *et al* (2022). We further note, that the time scale of decoupling is on the same order as the one reported for the silent alignment effect $t_{\text{sa}} = \frac{1}{s}$.

E.4. On-diagonal dynamics and the effect of initialisation variance

In this section we revisit the impact of initialisation scale for the on-diagonal dynamics. We now start with

$$a_1(t) = \frac{\frac{e^{-\frac{2s_2t}{\tau}} a_1(0)}{a_1(0)a_2(0)-b(0)^2} + \frac{1}{s_2} - \frac{1}{s_2} e^{-\frac{2s_2t}{\tau}}}{\frac{1}{a_1(0)a_2(0)-b(0)^2} \left(e^{-\frac{2(s_1+s_2)t}{\tau}} \left(1 - \frac{a_1(0)}{s_1} - \frac{a_2(0)}{s_2}\right) + e^{-\frac{2s_2t}{\tau}} \frac{a_1(0)}{s_1} + e^{-\frac{2s_1t}{\tau}} \frac{a_2(0)}{s_2} \right) + \frac{1}{s_2s_1}}}. \tag{200}$$

The diagonal elements simplify in the cases where $s_1 = s_2 = s$ (i.e. target \mathbf{Y} is diagonal),

$$a_1(t) = \frac{\frac{e^{-\frac{2st}{\tau}} a_1(0)}{a_1(0)a_2(0)-b(0)^2} + \frac{1}{s} - \frac{1}{s} e^{-\frac{2st}{\tau}}}{\frac{1}{a_1(0)a_2(0)-b(0)^2} \left(e^{-\frac{4st}{\tau}} \left(1 - \frac{a_1(0)}{s} - \frac{a_2(0)}{s}\right) + e^{-\frac{2st}{\tau}} \frac{a_1(0)}{s} + e^{-\frac{2st}{\tau}} \frac{a_2(0)}{s} \right) + \frac{1}{s^2}}}. \tag{201}$$

We consider when $|a_1(0)|, |a_2(0)|, |b(0)| \ll 1$, and recover a sigmoidal trajectory,

$$a_1(t) = \frac{sa_1(0)}{e^{-\frac{2st}{\tau}} [s - a_1(0) - a_2(0)] + a_1(0) + a_2(0)}. \tag{202}$$

Exact learning dynamics of deep linear networks with prior knowledge

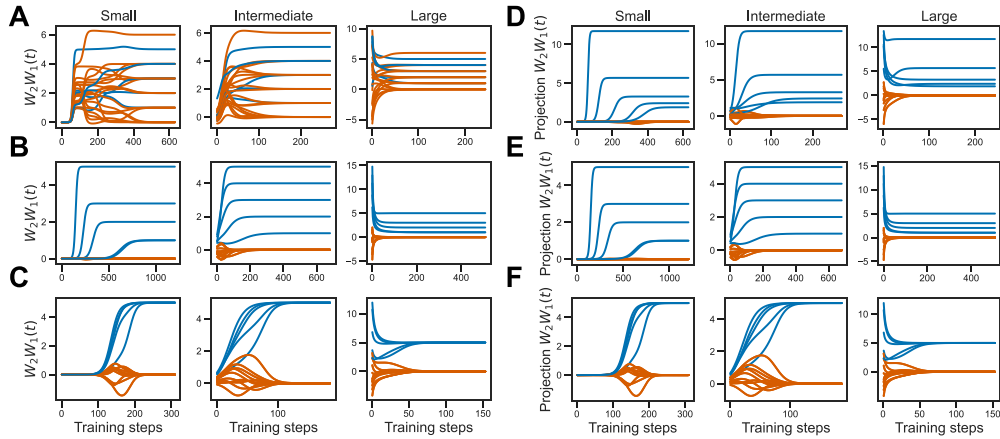


Figure 9. (A)–(C) Network function dynamics (Diagonal elements: blue, Off-diagonal elements: red) learning with learning rate $\eta = 0.01$ on the target 5×5 diagonal matrices shown in equation (198). The network was initialised as defined in section E with Small ($\sigma = 1 \times 10^{-6}$), Intermediate ($\sigma = 0.1$) and Large ($\sigma = 2$) variance, and hidden layer size $N_h = 10$. (A), Dense. (B), Diagonal. (C), Equal diagonal. (D)–(F). Corresponding numerical temporal dynamics of the projection of the network function on- and off-diagonal elements into the singular-basis of the initialisation. Equivalently, the temporal dynamics of the elements of $\mathbf{A}\mathbf{A}^T$ bottom left quadrant. (D), Dense. (E), Diagonal. (F), Equal diagonal.

We can compute the time at which $a_1(t)$ rises to half its asymptotic value to be

$$t_{\text{half}} = \frac{\tau}{2s} \log \left(\frac{s - a_1(0) - a_2(0)}{a_1(0) - a_2(0)} \right). \quad (203)$$

For $|a_1(0)|, |a_2(0)|, |b(0)| \gg 0$ the dynamics of the on-diagonal element a_1 is close to exponential.

The observation for 2×2 network hold for larger target matrices as shown in figure 9. For large variance initialisations, the dynamics are exponential. At intermediate variance initialisations, we observe more complex behaviour. While at small variance initialisations, the on-diagonal element describes a sigmoidal trajectory.

Appendix F. Continual learning

We consider the case of training a two-layer deep linear network on a sequence of tasks $\mathcal{T}_a, \mathcal{T}_b, \mathcal{T}_c, \dots$ with corresponding correlation functions $\mathcal{T}_a = \tilde{\Sigma}_a^{yx}, \mathcal{T}_b = \tilde{\Sigma}_b^{yx} \dots$. Then, the full batch loss of the i th task at any point in training time is

$$\mathcal{L}_i = \frac{1}{2P} \|\mathbf{W}_2 \mathbf{W}_1 \mathbf{X}_i - \mathbf{Y}_i\|_F^2. \quad (204)$$

From theorem 3.2 it follows that after training the network to convergence on task \mathcal{T}_j , the network function is $\mathbf{W}_2 \mathbf{W}_1 = \tilde{\mathbf{U}} \tilde{\mathbf{S}} \tilde{\mathbf{V}}^T = \tilde{\Sigma}_j^{yx}$. Further, using the assumption of whitened

inputs 2.2 and the identities $\|A\|_F^2 = \text{Tr}(AA^T)$ and $\text{Tr}(A) + \text{Tr}(B) = \text{Tr}(A + B)$, the full batch loss of the i -th task is then

$$\mathcal{L}_i(\mathcal{T}_j) = \frac{1}{2P} \left\| \tilde{\Sigma}_j^{yx} \mathbf{X}_i - \mathbf{Y}_i \right\|_F^2 \tag{205}$$

$$= \frac{1}{2P} \text{Tr} \left(\left(\tilde{\Sigma}_j^{yx} \mathbf{X}_i - \mathbf{Y}_i \right) \left(\tilde{\Sigma}_j^{yx} \mathbf{X}_i - \mathbf{Y}_i \right)^T \right) \tag{206}$$

$$= \frac{1}{2P} \text{Tr} \left(\tilde{\Sigma}_j^{yx} \mathbf{X}_i \mathbf{X}_i^T \tilde{\Sigma}_j^{yx^T} \right) - \frac{1}{P} \text{Tr} \left(\tilde{\Sigma}_j^{yx} \mathbf{X}_i \mathbf{Y}_i^T \right) + \frac{1}{2P} \text{Tr} \left(\mathbf{Y}_i \mathbf{Y}_i^T \right) \tag{207}$$

$$= \frac{1}{2} \text{Tr} \left(\tilde{\Sigma}_j^{yx} \tilde{\Sigma}_j^{yx^T} \right) - \text{Tr} \left(\tilde{\Sigma}_j^{yx} \tilde{\Sigma}_i^{yx^T} \right) + \frac{1}{2} \text{Tr} \left(\tilde{\Sigma}_i^{yy} \right) \tag{208}$$

$$= \frac{1}{2} \text{Tr} \left(\left(\tilde{\Sigma}_j^{yx} - \tilde{\Sigma}_i^{yx} \right) \left(\tilde{\Sigma}_j^{yx} - \tilde{\Sigma}_i^{yx} \right)^T - \tilde{\Sigma}_i^{yx} \tilde{\Sigma}_i^{yx^T} \right) + \frac{1}{2} \left(\tilde{\Sigma}_i^{yy} \right) \tag{209}$$

$$= \frac{1}{2} \left\| \tilde{\Sigma}_j^{yx} - \tilde{\Sigma}_i^{yx} \right\|_F^2 - \underbrace{\frac{1}{2} \text{Tr} \left(\tilde{\Sigma}_i^{yx} \tilde{\Sigma}_i^{yx^T} \right)}_c + \frac{1}{2} \left(\tilde{\Sigma}_i^{yy} \right). \tag{210}$$

Therefore, the amount of forgetting \mathcal{F} on task \mathcal{T}_i when training on task \mathcal{T}_k after having trained the network on task \mathcal{T}_j , i.e. the relative change of loss, is fully determined by the similarity structure of the tasks

$$\mathcal{F}_i(\mathcal{T}_j, \mathcal{T}_k) = \mathcal{L}_i(\mathcal{T}_k) - \mathcal{L}_i(\mathcal{T}_j) \tag{211}$$

$$= \frac{1}{2} \left\| \tilde{\Sigma}_k^{yx} - \tilde{\Sigma}_i^{yx} \right\|_F^2 + c - \frac{1}{2} \left\| \tilde{\Sigma}_j^{yx} - \tilde{\Sigma}_i^{yx} \right\|_F^2 - c \tag{212}$$

$$= \frac{1}{2} \left(\left\| \tilde{\Sigma}_k^{yx} - \tilde{\Sigma}_i^{yx} \right\|_F^2 - \left\| \tilde{\Sigma}_j^{yx} - \tilde{\Sigma}_i^{yx} \right\|_F^2 \right). \tag{213}$$

Appendix G. Revising structured knowledge

G.1. Reversal learning dynamics

In the following, we assume that the input dimension is equal to the output dimension. Further, we denote the i -th column of the left and right singular vectors as \mathbf{u}_i , $\tilde{\mathbf{u}}_i$ and \mathbf{v}_i , $\tilde{\mathbf{v}}_i$ respectively.

Reversal learning occurs when the task and the initial network function share the same left and right singular vectors, i.e. $\mathbf{U} = \tilde{\mathbf{U}}$ and $\mathbf{V} = \tilde{\mathbf{V}}$, except for one or multiple columns of the left singular vectors, for which the direction is reversed:

$$-\mathbf{u}_i = \tilde{\mathbf{u}}_i. \tag{214}$$

We note that, if there is any reversal in the right singular vectors $-\mathbf{v}_i = \tilde{\mathbf{v}}_i$, this can be written as a reversal in the left singular vectors, as the signs of the right and left singular

vectors are interchangeable. In the reversal learning setting, both $\mathbf{B} = \mathbf{U}^T \tilde{\mathbf{U}} + \mathbf{V}^T \tilde{\mathbf{V}}$ and $\mathbf{C} = \mathbf{U}^T \tilde{\mathbf{U}} - \mathbf{V}^T \tilde{\mathbf{V}}$ are diagonal matrices. The diagonal entries of \mathbf{C} are zero if the singular vectors are aligned and 2 if they are reversed. Similarly, diagonal entries of \mathbf{B} are 2 if the singular vectors are aligned and zero if they are reversed. Therefore, in the case of reversal learning, \mathbf{B} is a diagonal matrix with 0 values and thus is not invertible. As a consequence, the learning dynamics cannot be described by equation (11). 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 $\tilde{\mathbf{s}}_i$ denote the i -th diagonal entry of \mathbf{B} , \mathbf{C} , \mathbf{S} and $\tilde{\mathbf{S}}$ respectively, then the network dynamics can be rewritten as

$$\mathbf{W}_2 \mathbf{W}_1(t) = \frac{1}{2} \tilde{\mathbf{U}} \left(e^{\tilde{\mathbf{S}} \frac{t}{\tau}} \mathbf{B}^T + e^{-\tilde{\mathbf{S}} \frac{t}{\tau}} \mathbf{C}^T \right) \times \left[\mathbf{S}^{-1} + \frac{1}{4} \mathbf{B} \left(e^{2\tilde{\mathbf{S}} \frac{t}{\tau}} - \mathbf{I} \right) \tilde{\mathbf{S}}^{-1} \mathbf{B}^T - \frac{1}{4} \mathbf{C} \left(e^{-2\tilde{\mathbf{S}} \frac{t}{\tau}} - \mathbf{I} \right) \tilde{\mathbf{S}}^{-1} \mathbf{C}^T \right]^{-1} \tag{215}$$

$$\frac{1}{2} \left(e^{\tilde{\mathbf{S}} \frac{t}{\tau}} \mathbf{B} - e^{-\tilde{\mathbf{S}} \frac{t}{\tau}} \mathbf{C} \right) \tilde{\mathbf{V}}^T = \sum_{i=1}^{N_i} \frac{\mathbf{b}_i^2 e^{2\tilde{\mathbf{s}}_i \frac{t}{\tau}} - \mathbf{c}_i^2 e^{-2\tilde{\mathbf{s}}_i \frac{t}{\tau}}}{4\mathbf{s}_i^{-1} + \mathbf{b}_i^2 e^{2\tilde{\mathbf{s}}_i \frac{t}{\tau}} \tilde{\mathbf{s}}_i^{-1} - \mathbf{b}_i^2 \tilde{\mathbf{s}}_i^{-1} - \mathbf{c}_i^2 e^{-2\tilde{\mathbf{s}}_i \frac{t}{\tau}} \tilde{\mathbf{s}}_i^{-1} + \mathbf{c}_i^2 \tilde{\mathbf{s}}_i^{-1}} \tilde{\mathbf{u}}_i \tilde{\mathbf{v}}_i^T \tag{216}$$

$$= \sum_{i=1}^{N_i} \frac{\mathbf{s}_i \mathbf{b}_i^2 \tilde{\mathbf{s}}_i - \mathbf{s}_i \mathbf{c}_i^2 \tilde{\mathbf{s}}_i e^{-4\tilde{\mathbf{s}}_i \frac{t}{\tau}}}{4\tilde{\mathbf{s}}_i e^{-2\tilde{\mathbf{s}}_i \frac{t}{\tau}} + \mathbf{s}_i \mathbf{b}_i^2 \left(1 - e^{-2\tilde{\mathbf{s}}_i \frac{t}{\tau}} \right) + \mathbf{s}_i \mathbf{c}_i^2 \left(e^{-2\tilde{\mathbf{s}}_i \frac{t}{\tau}} - e^{-4\tilde{\mathbf{s}}_i \frac{t}{\tau}} \right)} \tilde{\mathbf{u}}_i \tilde{\mathbf{v}}_i^T. \tag{217}$$

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

$$\lim_{t \rightarrow \infty} \frac{-\mathbf{s}_i \mathbf{c}_i^2 \tilde{\mathbf{s}}_i e^{-4\tilde{\mathbf{s}}_i \frac{t}{\tau}}}{4\tilde{\mathbf{s}}_i e^{-2\tilde{\mathbf{s}}_i \frac{t}{\tau}} + \mathbf{s}_i \mathbf{c}_i^2 \left(e^{-2\tilde{\mathbf{s}}_i \frac{t}{\tau}} - e^{-4\tilde{\mathbf{s}}_i \frac{t}{\tau}} \right)} \tilde{\mathbf{u}}_i \tilde{\mathbf{v}}_i^T = 0. \tag{218}$$

Analytically, the learning dynamics are initialised and remain on the separatrix of a saddle point, until the corresponding singular value of the network function has vanished and remains zero, corresponding to convergence to the saddle point. When simulated numerically, the learning dynamics escape the saddle points due to imprecision of floating point arithmetic. However, numerical optimisation still suffers from catastrophic slowing (Lee *et al* 2022), as escaping the saddle point takes time (figure 6(A)). In contrast, in the case of aligned singular vectors ($\mathbf{c} = 0$), we recover the equation for the temporal dynamics as described in Saxe *et al* (2014). Training succeeds, as the singular value of the network function converges to its target value

$$\lim_{t \rightarrow \infty} \sum_{i=1}^{N_i} \frac{\mathbf{s}_i \mathbf{b}_i^2 \tilde{\mathbf{s}}_i}{4\tilde{\mathbf{s}}_i e^{-2\tilde{\mathbf{s}}_i \frac{t}{\tau}} + \mathbf{s}_i \mathbf{b}_i^2 \left(1 - e^{-2\tilde{\mathbf{s}}_i \frac{t}{\tau}} \right)} \tilde{\mathbf{u}}_i \tilde{\mathbf{v}}_i^T = \frac{\mathbf{s}_i \mathbf{b}_i^2 \tilde{\mathbf{s}}_i}{\mathbf{s}_i \mathbf{b}_i^2} \tilde{\mathbf{u}}_i \tilde{\mathbf{v}}_i^T \tag{219}$$

$$= \tilde{\mathbf{s}}_i \tilde{\mathbf{u}}_i \tilde{\mathbf{v}}_i^T. \tag{220}$$

In summary, in the case of aligned singular vectors, the learning dynamics can be described by the convergence of singular values. However in the case of reversal learning, analytically, training does not succeed. In simulations, the learning dynamics escape the saddle point due to numerical imprecision, but the learning dynamics are catastrophically slowed in the vicinity of the saddle point.

G.2. Exact learning dynamics in shallow networks

To provide a point of comparison to our deep linear network results, here we derive a solution for the temporal dynamics of reversal learning in a shallow network.

The network’s weights are optimised using full batch gradient descent with learning rate η (or equivalently time constant $\tau = 1/\eta$) on the mean squared error loss given in equation (2), yielding the first task dynamics

$$\tau \frac{d}{dt} \mathbf{W} = \tilde{\Sigma}^{yx} - \mathbf{W} \tilde{\Sigma}^{xx}, \tag{221}$$

where $\tilde{\Sigma}^{xx}$ and $\tilde{\Sigma}^{yx}$ is the input and input-output correlation matrices of the dataset. We define

$$\text{SVD}(\mathbf{W}(0)) = \mathbf{U}\mathbf{S}\mathbf{V}^T \text{ and } \text{SVD}(\tilde{\Sigma}^{yx}) = \tilde{\mathbf{U}}\tilde{\mathbf{S}}\tilde{\mathbf{V}}^T. \tag{222}$$

motivating the change of variable $\mathbf{W} = \mathbf{U}\overline{\mathbf{W}}\mathbf{V}^T$. We project the weight into the basis of the initialisation

$$\tau \frac{d}{dt} \mathbf{U}\overline{\mathbf{W}}\mathbf{V}^T = \tilde{\Sigma}^{yx} - \mathbf{U}\overline{\mathbf{W}}\mathbf{V}^T \tilde{\Sigma}^{xx} \tag{223}$$

$$\tau \frac{d}{dt} \mathbf{U}\overline{\mathbf{W}}\mathbf{V}^T = \mathbf{U}\mathbf{U}^T \tilde{\Sigma}^{yx} \mathbf{V}\mathbf{V}^T - \mathbf{U}\overline{\mathbf{W}}\mathbf{V}^T \tilde{\Sigma}^{xx} \tag{224}$$

$$\tau \frac{d}{dt} \overline{\mathbf{W}} = \mathbf{U}^T \tilde{\Sigma}^{yx} \mathbf{V} - \overline{\mathbf{W}} \tilde{\Sigma}^{xx}. \tag{225}$$

Under the assumption of whitened inputs 2.2, the dynamics yields

$$\tau \frac{d}{dt} \overline{\mathbf{W}} = \mathbf{U}^T \tilde{\Sigma}^{yx} \mathbf{V} - \overline{\mathbf{W}}. \tag{226}$$

Defining $\overline{\mathbf{W}}_{ii} = b_i$ the diagonal element of the matrix, encoding the strength of the mode i transmitted by the input-to-output weight. Similarly, we write $(\mathbf{U}^T \tilde{\Sigma}^{yx} \mathbf{V})_{ii} = k_i$. Assuming decoupled initial conditions, we obtain the scalar dynamics

$$\tau \frac{d}{dt} b_i = k_i - b_i \tag{227}$$

with solution

$$b_i = k_i \left(1 - e^{-\frac{t}{\tau}} \right) + b_i^0 e^{-\frac{t}{\tau}}. \tag{228}$$

Reverting the change of variable, the weight trajectory yields

$$\mathbf{W} = \mathbf{U}\mathbf{B}(t)\mathbf{V}^T. \quad (229)$$

This solution is very similar to the one proposed by Saxe *et al* (2019). However, the key here is that k_i can have negative values. k_i is negative whenever a vector is in the opposite direction to the initialisation (as in the reversal learning setting). We show in figure 6 that the analytical solution derived above matches the numerical temporal dynamics. From equation (228), we note that the shallow network cannot display catastrophic slowing.

Appendix H. Simulations

In the following, we describe the details of the simulation studies. Generally, N_i , N_h and N_o denote the dimension of the input, hidden layer and output (target) respectively. The number of training samples is N and the learning rate is denoted by $\eta = \frac{1}{\tau}$.

H.1. Zero-balanced weight initialisation

The initial network weights are zero-balanced 2.3 when they satisfy

$$\mathbf{W}_1(0)\mathbf{W}_1(0)^T = \mathbf{W}_2(0)^T\mathbf{W}_2(0). \quad (230)$$

In practice, we use algorithm 1 to initialise the network weights, where α is a scaling factor which is used to control the variance of the weights, i.e. to vary between small and large weight initialisations.

Algorithm 1. Zero-balanced weight initialisation.

```

Require:  $N_i, N_h, N_o, \sigma$ 
 $\mathbf{W}_1 \sim \mathcal{N}(\mu = 0, \sigma) \in \mathbb{R}^{N_h \times N_i}$ 
 $\mathbf{W}_2 \sim \mathcal{N}(\mu = 0, \sigma) \in \mathbb{R}^{N_o \times N_h}$ 
 $\mathbf{U}, \mathbf{S}, \mathbf{V} \leftarrow \text{SVD}(\mathbf{W}_2\mathbf{W}_1)$ 
 $\mathbf{S} \leftarrow \sqrt{\mathbf{S}}$ 
 $\mathbf{R} \sim \mathcal{N}(\mu = 0, \sigma = 1) \in \mathbb{R}^{N_h \times N_h}$ 
 $\mathbf{R}_{\cdot, -} \leftarrow \text{SVD}(\mathbf{R})$ 
if  $N_i \neq N_o$  then
   $N_s \leftarrow N_i$  if  $N_i < N_o$  else  $N_o$ 
   $\mathbf{S}_1 \leftarrow \begin{bmatrix} \mathbf{S} \\ \mathbf{0}_{N_h - N_s \times N_s} \end{bmatrix}$ 
   $\mathbf{S}_2 \leftarrow \begin{bmatrix} \mathbf{S} & \mathbf{0}_{N_s \times N_h - N_s} \end{bmatrix}$ 
   $\mathbf{W}_1 \leftarrow \mathbf{R}\mathbf{S}_1\mathbf{V}^T$ 
   $\mathbf{W}_2 \leftarrow \mathbf{U}\mathbf{S}_2\mathbf{R}^T$ 
else
   $\mathbf{W}_1 \leftarrow \mathbf{R}\mathbf{S}\mathbf{V}^T$ 
   $\mathbf{W}_2 \leftarrow \mathbf{U}\mathbf{R}^T$ 
end if
return  $\mathbf{W}_1\mathbf{W}_2$ 

```

H.2. Tasks

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

H.2.1. Random regression task. In a random regression task the inputs $\mathbf{X} \in \mathbb{R}^{N_i, N}$ are sampled from a random normal distribution $\mathbf{X} \sim \mathcal{N}(\mu = 0, \sigma = 1)$. The input data \mathbf{X} is then whitened, such that $\frac{1}{N}\mathbf{X}\mathbf{X}^T = \mathbf{I}$. The target values $\mathbf{Y} \in \mathbb{R}^{N_o, N}$ are also sampled from a random normal distribution, however, with variance adjusted to the number of output nodes $\mathbf{Y} \sim \mathcal{N}(\mu = 0, \alpha = \frac{1}{\sqrt{N_o}})$. Thus, network inputs and target values are uncorrelated Gaussian noise and therefore, a linear solution does not always exist.

H.2.2. Teacher-student task. In order to guarantee that a linear solution exists, we use the teacher-student setup. First, inputs \mathbf{X} are sampled as in the random regression task. Then, target values \mathbf{Y} are generated by sampling a pair of random zero-balanced weights $\mathbf{W}_1 \in \mathbb{R}^{N_h \times N_i}$ and $\mathbf{W}_2 \in \mathbb{R}^{N_o \times N_h}$ and then calculating $\mathbf{Y} = \mathbf{W}_2\mathbf{W}_1\mathbf{X}$. Like this, it is ensured that a linear solution exists. The variance of the output is varied by changing the variation within the zero-balanced weights σ .

H.2.3. Semantic hierarchy. Input items in the semantic hierarchy task are encoded as one-hot vectors, i.e. $\mathbf{X} = \mathbf{I}$. The corresponding target vectors y_i encoded the position in the hierarchical tree. Where a 1 encoded being a left child of a node, a -1 encoded being a right child of a node and a 0 encoded that the item is not a child of that node. For example, the blue fish is a blue fish, it is a left child of the root node, a left child of the animal node, not part of the plant branch, a right child of the fish node, and not part of the bird, algae or flower branch, leading to the label $[1, 1, 1, 0, -1, 0, 0, 0]$. The labels for all objects in the semantic tree as depicted in figure 3(A) is then

$$\mathbf{Y} = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 & -1 & -1 & -1 & -1 \\ 1 & 1 & -1 & -1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 1 & -1 & -1 \\ 1 & -1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & -1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & -1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & -1 \end{bmatrix}. \tag{231}$$

The singular value decomposition for the corresponding correlation matrix $\tilde{\Sigma}^{yx}$ are not unique. The first two, the third and the fourth and the last four singular values are identical. In order to match the numerical and analytical solution, this permutation invariance is removed by adding a small constant perturbation to each column $\mathbf{y}_i, i \in 1, \dots, N$ of the labels

$$\mathbf{y}_i = \mathbf{y}_i * \left(1 + \frac{0.1}{i}\right), \tag{232}$$

leading to almost but not exactly identical singular values.

H.2.4. Colour hierarchy. Following the same procedure as described for the semantic hierarchy, the labels for the colour hierarchy as depicted in figure 6(C) are then

$$\mathbf{Y} = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\ -1 & 1 & 1 & -1 & 1 & 1 & -1 & -1 \\ 0 & -1 & 1 & 0 & -1 & 1 & 0 & 0 \\ 1 & 0 & 0 & -1 & 0 & 0 & -1 & 1 \\ 0 & 0 & 1 & 0 & 0 & -1 & 0 & 0 \\ 0 & -1 & 0 & 0 & 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 & -1 \\ 0 & 0 & 0 & 1 & 0 & 0 & -1 & 0 \end{bmatrix}. \quad (233)$$

H.3. Figure 1

Figure 1 panels (B)–(D) show three simulations from varying initial weights on the same teacher-student task. The task was created with $\sigma = 0.35$. Further, $N_i = 5$, $N_h = 10$, $N_o = 2$ and $N = 10$. The learning rate was $\eta = 0.1$ and the initial network weights were sampled with $\sigma = 0.01$, $\sigma = 0.25$ and $\sigma = 0.25$ in panels (B), (C) and (D) respectively.

H.4. Figure 2

Figure 2 panels (A) and (B) show a simulation on the same teacher-student task ($\sigma = 0.25$), once from small initial weights ($\sigma = 0.01$) and once from large initial weights ($\sigma = 0.15$). Dimensions were $N_i = 4$, $N_h = 5$, $N_o = 3$ and $N = 10$ and the learning rate was $\eta = 0.05$. Panel (C) was generated by running 50 simulations, each with a different initial random seed. For each of the simulations, dimensions were sampled randomly, such that $N_i \in [2, 50]$, $N_o \in [2, 50]$, $N_h = [\min(N_i, N_o), 50]$ and $N \in [2 \max(N_i, N_h, N_o), 3 \max(N_i, N_h, N_o)]$. Then, a random regression task was generated. Subsequently, a linear network was initialised with $\sigma \sim \mathcal{U}[\frac{0.01}{\sqrt{\max(N_i, N_o, N_h)}}, \frac{0.5}{\sqrt{\max(N_i, N_o, N_h)}}]$. The network was then trained until convergence on the same task from the same initial weights for seven different learning rates $\eta \in \{0.05, 0.0232, 0.0107, 0.005, 0.0023, 0.0011, 0.0005\}$.

H.5. Figure 3

Panels (C)–(F) in figure 3 were generated by training a linear network with $N_i = 8$, $N_h = 14$, $N_o = 8$ on the $N = 8$ items of the semantic hierarchy task. The learning rate was $\eta = 0.05$ and the initial weights in panels (C), (D) and (E) were sampled from a normal distribution with $\sigma = 0.0001$ and $\sigma = 0.42$ and zero-balanced weights with $\sigma = 0.44$ respectively.

H.6. Figure 4

Figure 4 panel (A) was generated by training a linear network with $N_i = 5$, $N_h = 10$, $N_o = 5$ on the target \mathbf{Y} as shown in equation (198) (equal diagonal). The network was initialised with $\sigma = 0.1$. The learning rate was $\eta = 0.01$.

Figure 4 panels (D)–(F) was generated by training a linear network with $N_i = 2$, $N_h = 10$, $N_o = 2$ on the target \mathbf{Y} as shown in figure 4(C) and input $\mathbf{X} = bfi$. The network was initialised with small $\sigma = 0.000\ 01$, intermediate $\sigma = 0.3$ and large $\sigma = 2$ synaptic weights. The learning rate was $\eta = 0.0001$.

H.7. Figure 5

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

Panels (B) and (C) were generated by running 50 simulations on two subsequent random regression tasks, each with a different initial random seed. The simulation was repeated three times, the first time with a linear, the second time with a tanh and the last time with a ReLU activation function in the hidden layer. Dimension were randomly sampled such that $N_i \in [2, 30]$, $N_o \in [2, 30]$, $N_h = [\min(N_i, N_o), 30]$ and $N = 100$. The standard deviation of the initial weight was chosen such that $\sigma = \frac{0.5}{\sqrt{0.5(N_i + N_h)}}$. The learning rate was $\eta = 0.075$.

For panel (D) and (E) the same simulation was repeated for three times, the first time with a linear, the second time with a tanh and the last time with a ReLU activation function. Each time, five random regression tasks with dimensions $N_i = 15$, $N_h = 18$, $N_o = 21$ and $N = 50$ were generated. Then a network with initial weight scale $\alpha = 0.025$ was sequentially trained with learning rate $\eta = 0.1$ on the five random regression tasks.

H.8. Figure 6

Figure 6 panel (A) was generated by training a linear network with $N_i = 4$, $N_h = 6$, $N_o = 4$ on a reversal learning task (see section G.1), which was derived from a random regression task. The learning rate was $\eta = 0.05$ and initial weights had a standard deviation of $\sigma = 0.25$. Panel (B) was generated by training a shallow linear network (see section G.2) on the same reversal learning task, with identical hyperparameters as in panel (A).

For the top and bottom rows of panels (E) and (F) a linear network with $N_i = 8$, $N_h = 14$, $N_o = 8$ was trained on the semantic hierarchy task, followed by training the network on the adapted semantic hierarchy as depicted in figure 6(C) top, which is a reversal learning task and the colour hierarchy respectively. The learning rate was $\eta = 0.05$ and σ was set to 0.001 and 0.35 respectively.

References

- Arora R *et al* 2020 (Princeton University) Theory of deep learning (in preparation)
 Arora S, Cohen N, Golowich N and Wei H 2018b A convergence analysis of gradient descent for deep linear neural networks (arXiv:1810.02281)
 Arora S, Cohen N and Hazan E 2018a On the optimization of deep networks: implicit acceleration by overparameterization *Int. Conf. on Machine Learning* (PMLR) pp 244–53
 Arora S, Cohen N, Wei H and Luo Y 2019a Implicit regularization in deep matrix factorization *Advances in Neural Information Processing Systems* vol 32
 Arora S, Du S S, Wei H, Zhiyuan Li, Salakhutdinov R R and Wang R 2019b On exact computation with an infinitely wide neural net *Advances in Neural Information Processing Systems* vol 32

- Asanuma H, Takagi S, Nagano Y, Yoshida Y, Igarashi Y and Okada M 2021 Statistical mechanical analysis of catastrophic forgetting in continual learning with teacher and student networks *J. Phys. Soc. Japan* **90** 104001
- Atanasov A, Bordelon B and Pehlevan C 2022 Neural networks as kernel learners: the silent alignment effect *Int. Conf. on Learning Representations*
- Bahri Y, Kadmon J, Pennington J, Schoenholz S S, Sohl-Dickstein J and Ganguli S 2020 Statistical mechanics of deep learning *Annu. Rev. Condens. Matter Phys.* **11** 501–28
- Baldi P and Hornik K 1989 Neural networks and principal component analysis: learning from examples without local minima *Neural Netw.* **2** 53–58
- Bengio Y, Louradour Jôme, Collobert R and Weston J 2009 Curriculum learning *Proc. 26th Annual Int. Conf. on Machine Learning* pp 41–48
- Biehl M and Schwarze H 1995 Learning by on-line gradient descent *J. Phys. A: Math. Gen.* **28** 643
- Carey S E 1985 *Conceptual Change In Childhood* (MIT Press)
- Carleo G, Cirac I, Cranmer K, Daudet L, Schuld M, Tishby N, Vogt-Maranto L and Zdeborová L 2019 Machine learning and the physical sciences *Rev. Mod. Phys.* **91** 045002
- Chizat L, Oyallon E and Bach F 2019 On lazy training in differentiable programming *Advances in Neural Information Processing Systems* vol 32
- Doan T, Abbana Bennani M, Mazoure B, Rabusseau G and Alquier P 2021 A theoretical analysis of catastrophic forgetting through the ntk overlap matrix *Int. Conf. on Artificial Intelligence and Statistics* (PMLR) pp 1072–80
- Erdeniz B and Bedin Atalay N 2010 Simulating probability learning and probabilistic reversal learning using the attention-gated reinforcement learning (agrel) model *2010 Int. Joint Conf. on Neural Networks (IJCNN)* (IEEE) pp 1–6
- Flesch T, Balaguer J, Dekker R, Nili H and Summerfield C 2018 Comparing continual task learning in minds and machines *Proc. Natl Acad. Sci.* **115** E10313–22
- Flesch T, Juechems K, Dumbalska T, Saxe A and Summerfield C 2022 Orthogonal representations for robust context-dependent task performance in brains and neural networks *Neuron* **110** 4212–19
- French R M 1999 Catastrophic forgetting in connectionist networks *Trends Cogn. Sci.* **3** 128–35
- Fukumizu K 1998 Effect of batch learning in multilayer neural networks *Int. Conf. Neural Information Processing (ICONIP)* pp 67–70
- Gerace F, Saglietti L, Sarao Mannelli S, Saxe A and Zdeborová L 2022 Probing transfer learning with a model of synthetic correlated datasets *Mach. Learn.: Sci. Technol.* **3** 015030
- Glorot X and Bengio Y 2010 Understanding the difficulty of training deep feedforward neural networks *Proc. Thirteenth 13th Int. Conf. on Artificial Intelligence and Statistics (JMLR Workshop and Conf. Proc.)* pp 249–56
- Goldt S, Advani M, Saxe A M, Krzakala F and Zdeborová L 2019 Dynamics of stochastic gradient descent for two-layer neural networks in the teacher-student setup *Advances in Neural Information Processing Systems* vol 32
- Gunasekar S, Lee J D, Soudry D and Srebro N 2018 Implicit bias of gradient descent on linear convolutional networks *Advances in Neural Information Processing Systems* p 31
- Huh D 2020 Curvature-corrected learning dynamics in deep neural networks *Int. Conf. on Machine Learning* (PMLR) pp 4552–60
- Jacot A, Gabriel F and Hongler C 2018 Neural tangent kernel: Convergence and generalization in neural networks *Advances in Neural Information Processing Systems* vol 31
- Javed K and White M 2019 Meta-learning representations for continual learning *Advances in Neural Information Processing Systems* pp 1820–30
- Kaiming H, Zhang X, Ren S and Sun J 2015 Delving deep into rectifiers: Surpassing human-level performance on imagenet classification *Proc. IEEE Int. Conf. on Computer Vision* pp 1026–34
- Kirkpatrick J *et al* 2017 Overcoming catastrophic forgetting in neural networks *Proc. Natl Acad. Sci.* **114** 3521–6
- Kriegeskorte N, Mur M and Bandettini P A 2008 Representational similarity analysis-connecting the branches of systems neuroscience *Front. Syst. Neurosci.* **2** 4
- Lampinen A K and Ganguli S 2018 An analytic theory of generalization dynamics and transfer learning in deep linear networks (arXiv:1809.10374)
- Laurent T and Brecht J 2018 Deep linear networks with arbitrary loss: all local minima are global *Int. Conf. on Machine Learning* (PMLR) pp 2902–7
- Lee J, Xiao L, Schoenholz S, Bahri Y, Novak R, Sohl-Dickstein J and Pennington J 2019 Wide neural networks of any depth evolve as linear models under gradient descent *Advances in Neural Information Processing Systems* vol 32
- Lee S, Sarao Mannelli S, Clopath C, Goldt S, and Saxe A 2022 Maslow’s hammer for catastrophic forgetting: node re-use vs node activation (arXiv:2205.09029)
- Lee S, Sebastian G and Saxe A 2021 Continual learning in the teacher-student setup: Impact of task similarity *Int. Conf. on Machine Learning* (PMLR) pp 6109–19
- McClelland J L 2013 Incorporating rapid neocortical learning of new schema-consistent information into complementary learning systems theory *J. Exp. Psychol. Gen.* **142** 1190

- McClelland J L, McNaughton B L and O'Reilly R C 1995 Why there are complementary learning systems in the hippocampus and neocortex: insights from the successes and failures of connectionist models of learning and memory *Psychol. Rev.* **102** 419
- McCloskey M and Cohen N J 1989 Catastrophic interference in connectionist networks: The sequential learning problem *Psychology of Learning and Motivation* vol 24 (Elsevier) pp 109–65
- Mei S, Montanari A and Nguyen P-M 2018 A mean field view of the landscape of two-layer neural networks *Proc. Natl Acad. Sci.* **115** E7665–71
- Mishkin D and Matas J 2015 All you need is a good init (arXiv:1511.06422)
- Murphy G 2004 *The big Book of Concepts* (MIT Press)
- Parisi G I, Kemker R, Part J L, Kanan C and Wermter S 2019 Continual lifelong learning with neural networks: a review *Neural Netw.* **113** 54–71
- Pennington J, Schoenholz S and Ganguli S 2017 Resurrecting the sigmoid in deep learning through dynamical isometry: theory and practice *Advances in Neural Information Processing Systems* vol 30
- Poggio T, Liao Q, Miranda B, Banburski A, Boix X and Hidary J 2018 Theory iiib: generalization in deep networks (arXiv:1806.11379)
- Raghu M, Zhang C, Kleinberg J and Bengio S 2019 Transfusion: understanding transfer learning for medical imaging *Advances in Neural Information Processing Systems* p 32
- Ratcliff R 1990 Connectionist models of recognition memory: constraints imposed by learning and forgetting functions *Psychol. Rev.* **97** 285
- Rotskoff G and Vanden-Eijnden E 2018 Parameters as interacting particles: long time convergence and asymptotic error scaling of neural networks *Advances in Neural Information Processing Systems* vol 31
- Saad D and Solla S A 1995 Exact solution for on-line learning in multilayer neural networks *Phys. Rev. Lett.* **74** 4337
- Saxe A M, McClelland J L and Ganguli S 2014 Exact solutions to the nonlinear dynamics of learning in deep linear neural networks *2nd Int. Conf. on Learning Representations, ICLR 2014 (Conf. Track Proc.) (Banff, AB, Canada, 14–16 April 2014)*
- Saxe A M, McClelland J L and Ganguli S 2019 A mathematical theory of semantic development in deep neural networks *Proc. Natl Acad. Sci.* **116** 11537–46
- Shachaf G, Brutzkus A and Globerson A 2021 A theoretical analysis of fine-tuning with linear teachers *Advances in Neural Information Processing Systems* vol 34
- Simon D and Wei H 2019 Width provably matters in optimization for deep linear neural networks *Int. Conf. on Machine Learning* (PMLR) pp 1655–64
- Sirignano J and Spiliopoulos K 2020 Mean field analysis of neural networks: A central limit theorem *Stoch. Process. Appl.* **130** 1820–52
- Tarmoun S, Franca G, Haeffele B D and Vidal R 2021 Understanding the dynamics of gradient flow in overparameterized linear models *Int. Conf. on Machine Learning* (PMLR) pp 10153–61
- Taylor M E and Stone P 2009 Transfer learning for reinforcement learning domains: a survey *J. Mach. Learn. Res.* **10** 1633–85
- Thrun S and Pratt L 2012 *Learning to Learn* (Springer Science & Business Media)
- Tripuraneni N, Jordan M and Jin C 2020 On the theory of transfer learning: The importance of task diversity *Advances in Neural Information Processing Systems* vol 33 pp 7852–62
- Xiao L, Bahri Y, Sohl-Dickstein J, Schoenholz S and Pennington J 2018 Dynamical isometry and a mean field theory of cnns: How to train 10,000-layer vanilla convolutional neural networks *Int. Conf. on Machine Learning* (PMLR) pp 5393–402
- Yan W-Y, Helmke U and Moore J B 1994 Global analysis of oja's flow for neural networks *IEEE Trans. Neural Netw.* **5** 674–83
- Zenke F, Poole B and Ganguli S 2017 Continual learning through synaptic intelligence *Int. Conf. on Machine Learning* (PMLR) pp 3987–95
- Ziwei J and Telgarsky M 2018 Gradient descent aligns the layers of deep linear networks (arXiv:1810.02032)