
A Theory of Initialisation’s Impact on Specialisation

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Prior work has demonstrated a consistent tendency in neural networks engaged in
2 continual learning tasks, wherein intermediate task similarity results in the highest
3 levels of catastrophic interference with prior learning. This phenomenon is at-
4 tributed to the network’s tendency to reuse learned features across tasks. However,
5 this explanation heavily relies on the condition that such a neuron specialisation oc-
6 curs, i.e. the emergence of localised representations. Our investigation challenges
7 the validity of this assumption. Using theoretical frameworks for the analysis of neu-
8 ral networks, we show a strong dependence of specialisation on the initial condition.
9 More precisely, we show that weight imbalance and high weight entropy can favour
10 specialised solutions. We then apply these insights in the context of continual learn-
11 ing, first showing the emergence of a monotonic relation between task-similarity
12 and forgetting in non-specialised networks, and, finally, assessing the implications
13 on the commonly employed elastic weight consolidation regularisation technique.

14 1 Introduction

15 Theories of representation in biological neural networks span from highly localised representations in
16 single neural units [1] to fully distributed or shared representations [2]. While shared representations
17 offer greater resilience, specialised representations allow for more efficient encoding of information.
18 Experimental evidence supports both ends of this spectrum, with different brain areas and tasks
19 exhibiting distinct forms of representation [3, 4, 5, 6, 7]. Similarly, artificial neural networks display
20 both shared [8, 9, 10] and specialised representations [11, 12], where a recent advancements in
21 explainable AI, such as the Golden Gate Claude model [13], exemplify an extreme of the spectrum.

22 Given the trade-offs between shared and specialised representations, a critical research challenge
23 lies in understanding how to guide neural networks towards one form or the other. This tension is
24 especially relevant in contexts like disentangled representation learning [14] and multi-task learning
25 [15], including continual learning and transfer learning. Specialised representations can facilitate
26 faster adaptation and reduce catastrophic forgetting [16, 17], as they allow networks to rewire
27 efficiently [18]. Rich Caruana’s seminal work on multi-task learning [15] emphasised the value of
28 specialisation in enhancing performance across multiple tasks. Recent efforts to mitigate catastrophic
29 forgetting [19, 20] have led to the development of regularisation strategies that promote specialisation,
30 such as elastic weight consolidation [21], synaptic intelligence [22], and learning without forgetting
31 [23]. In disentangled representation learning, [24] highlighted that, despite the potential success of
32 unsupervised approaches, disentanglement does not emerge naturally without an explicit inductive
33 bias, underscoring the need for supervision to enforce such structures.

34 In this study, we investigate the role of initialisation in steering neural networks towards specialised
35 or shared representations, providing a complementary perspective on both the lazy learning regime
36 [25] and the rich learning regime [26, 27, 28]. Previous research [29, 30, 31] has shown that by
37 interpolating between these regimes, we can transition from shared representations—characterised by
38 random projections in the neural tangent kernels—to effective feature learning [32, 33, 34, 35, 36].

39 While our analysis remains within the feature learning regime, it adopts a distinct theoretical approach
40 compared these studies, concentrating specifically on the impact of initialisation within standard
41 synthetic frameworks for neural networks. This exploration reveals how initialisation can skew
42 the learning dynamics towards either specialised or shared representations, thereby adding a new
43 dimension to the study of learning dynamics in over-parameterised networks.

44 Our work makes the following **main contributions**:

- 45 • We study the impact of initialisation on specialisation through two theoretical frameworks:
 - 46 – We utilise the dynamics of **deep linear networks** to investigate the evolution of
47 specialisation [37];
 - 48 – We extend this analysis to **high-dimensional mean-field neural networks**, drawing
49 insights from stochastic gradient dynamics [38, 39, 40].
- 50 • Our findings challenge prevailing assumptions regarding the relationship between task
51 similarity and catastrophic forgetting [41, 42, 43].
- 52 • Moreover, we identify specific initialisation schemes that promote specialised solutions by
53 increasing the entropy of the readout weights and creating an imbalance between the first
54 and last layers, akin to the findings of [35].
- 55 • Finally, we demonstrate the practical implications of our results on regularisation strategies,
56 specifically analysing how Elastic Weight Consolidation (EWC) [21] is influenced by spe-
57 cialisation dynamics, highlighting potential pitfalls associated with regularisation methods
58 in continual learning.

59 In Sec. 2, we introduce the concept of specialisation within the teacher-student framework and
60 highlight the relevant literature. Sec. 3 explores this issue through the lens of deep linear dynam-
61 ics, illustrating its impact on learned representations, particularly in the context of disentangled
62 representation learning. Sec. 4 addresses the continual learning problem, revisiting existing theoret-
63 ical frameworks and demonstrating how their conclusions may not hold under certain initialisation
64 schemes. We conclude this section by discussing the implications for the EWC mitigation strategy.
65 Finally, in Sec. 5, we reflect on the limitations of our work and propose future directions for research.

66 2 Specialisation in the teacher-student

67 The teacher-student framework is a generative model that allows for the controlled creation of
68 synthetic datasets [44]. The framework involves two classifiers: the *teacher* and the *student*, for
69 instance represented as neural networks as exemplified in Fig. 1a. The teacher, has fixed randomly
70 drawn weights and maps random inputs \mathbf{x} from a given distribution to labels, providing a rule for
71 generating data. The student, on the other hand, updates its parameters through learning protocols
72 like stochastic gradient descent (SGD) to approximate the teacher’s outputs.

73 While a detailed quantitative characterisation of specialisation follows in the next sections, we briefly
74 introduce the concept within the teacher-student framework. [38] showed that, when both teacher
75 and student are modelled as committee machines, each student neuron specialises by aligning with a
76 specific teacher neuron. Similarly, [45] observed that for certain activation functions in two-layer
77 networks, an over-parameterised student will selectively use only a subset of those units to replicate
78 the teacher’s outputs. This phenomenon, termed specialisation, stands in contrast to a student
79 redundantly sharing representations of the teacher across neurons. In this work we present a more
80 comprehensive account of the factors underlying specialisation. In contrast to [45], we argue that
81 initialisation—not the activation function—is chiefly responsible. We highlight this in Fig. 1b, by
82 showing that with carefully chosen initialisations we can train a highly specialised ReLU student
83 (bottom panels), and a non-specialising sigmoidal student (top), which represents the opposite of the
84 conclusions presented in [45].

85 3 Specialisation explained using Linear Dynamics

86 As a first step towards understanding specialisation in neural networks we turn to the deep linear
87 neural network paradigm [37]. While deep linear networks can only represent linear input-output
88 mappings, they showcase intricate fixed point structure and nonlinear learning dynamics reminiscent

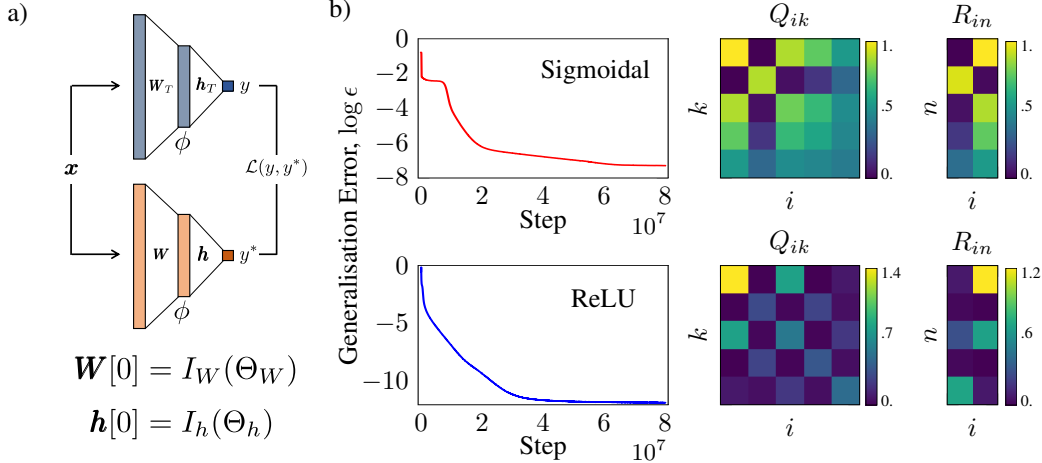


Figure 1: **Initialisation impacts specialisation.** a) In the teacher-student setup a student network is trained with labels generated by a fixed teacher network. Previous work established a relationship between the activation function ϕ and the propensity for the student nodes to specialise to teacher nodes. However we show in this work that this is an overly simplistic description; other factors including student weight initialisations I_W, I_h , parameterised by Θ_W, Θ_h arguably play a stronger role. b) Generalisation error curves for two simulations of the teacher-student setup, one with a ReLU activation function and one with a scaled error activation function. Θ_W and Θ_h are chosen to achieve a solution with ReLU that specialises—as indicated by sparser overlap matrices on the bottom right, and a scaled error function solution that does not specialise—as indicated by denser overlap matrices on the top right. A sparse (dense) Q matrix shows few (many) nodes are active, while a sparse (dense) R matrix shows student nodes are representing teacher nodes in a targeted (redundant) manner. Further details for the quantities described can be found in Sec. 4.

89 of phenomena seen in nonlinear networks. Deep linear networks have been successfully used to
 90 describe the effects of depth and nonlinearity, while showcasing the influence of initialisation [46, 47].
 91 Here we construct a synthetic setup, to study the influence of initialisation on specialisation. In this
 92 work, we consider specialisation adhering to the definition of proposed by the statistical physics
 93 literature [45] which considers whether one neuron will account for all of the variance associated to
 94 one feature, while the others remain inactive. This is in contrast to other work on modularity [48]
 95 such as Neural Module Networks [49, 50, 51, 52], mixture-of-expert models [53, 54, 55], tensor
 96 product networks [56], among others [57, 58], which consider specialisation as a subset of a network
 97 or module performing a single “task” or only being activated by one interpretable feature in the
 98 dataset. Thus, these works are more concerned with *what* is learned and consider specialisation to
 99 imply feature sparsity [59]. While we are concerned with the manner in which learning is represented,
 100 a phenomenon closer to activation sparsity.

101 3.1 Specialisation in the deep linear network framework

102 To connect this framework to specialisation we use the notion of the “neural race” from [46]. The
 103 neural race hypothesis says that the pathways through a network are racing to explain the variance in
 104 the dataset (perform the input-output mapping). Thus, we consider the limited case of a network with
 105 two hidden neurons and one output neuron. Fig. 2 depicts the setup, notation and strategy for this
 106 section. We ask the question: “when will one pathway finish learning (reach it’s hitting time t^*) before
 107 the other begins learning (reaches it’s escaping time \hat{t})”. In cases when this occurs, the network would
 108 have specialised as only one pathway will have any activity and will explain all of the data. Similar
 109 to Sec. 4 we generate data by sampling the elements of a data point from a Gaussian distribution
 110 ($x_i \sim \mathcal{N}(0, 1)$) with $i = 1, \dots, d$. We then define a ground-truth mapping (\mathbf{W}_T) and generate labels
 111 $y = \mathbf{W}_T \cdot \mathbf{x}$. We only consider regression tasks in this section, thus $y \in \mathbb{R}$. For P inputs we can
 112 form the input matrix $\mathbf{X} \in \mathbb{R}^{d \times P}$ and row vector of scalar outputs $\mathbf{y} \in \mathbb{R}^{1 \times P}$. The dataset statistics
 113 which drive learning are collected in the input and input-output correlation matrices, Σ^x and Σ^{yx}

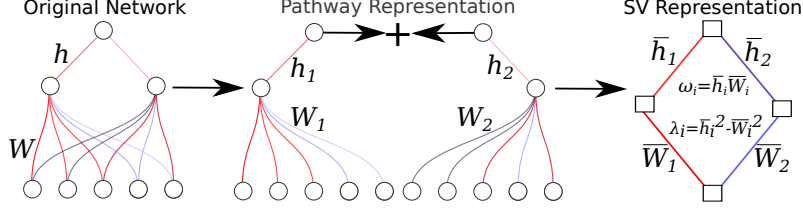


Figure 2: Summary of our setup, notation and strategy. a) The original network with two hidden neurons learning the regression task. b) We split the network into two separate pathways and consider their dynamics individually. Since both networks are learning the same task simultaneously, their dynamics are coupled. c) To obtain the dynamics of the two pathways and calculate their escaping and hitting time we track the pathway dynamics in terms of the network’s effective singular values. The closed form dynamics for the pathway singular value are given in Eq. 3.

114 respectively. For the task described above the singular value decomposition of these matrices are:

$$\Sigma^x = E[\mathbf{X}\mathbf{X}^T] = \mathbf{V}\mathbf{D}\mathbf{V}^T, \quad \Sigma^{yx} = E[\mathbf{y}\mathbf{X}^T] = \mathbf{u}\mathbf{s}\mathbf{v}^T. \quad (1)$$

115 Here, $\mathbf{u} \in \{-1, 1\}$, \mathbf{v} is a vector such that $\mathbf{v}^T\mathbf{v} = 1$ and \mathbf{V} is an orthogonal singular vector matrix.
 116 Correspondingly, s is the singular value for the rank 1 task and \mathbf{D} is a diagonal matrix of singular
 117 values. Note that we assume that the correlation matrices are mutually diagonalisable (share the
 118 same \mathbf{V}) up to the rank of Σ^{yx} .

119 For this task we consider a single hidden layer network (Fig. 2 left) computing output $\hat{y} = \mathbf{h}\mathbf{W}\mathbf{x}$ with
 120 $\mathbf{h} \in \mathbb{R}^K$ and $\mathbf{W} \in \mathbb{R}^{K \times d}$ in response to an input $\mathbf{x} \in \mathbb{R}^d$. The network is trained to minimise the
 121 mean squared error loss using full batch gradient descent with a small learning rate η . To identify
 122 when specialisation will occur in this network, we split the network into two pathways with one
 123 hidden neuron each. The input and output dimensions remain the same (Fig. 2 middle). Finally
 124 we obtain the linear dynamics (ultimately depicted as Eq. 3) for each pathway (the full details and
 125 assumptions of the derivation are given in Appendix A). In this setting, the network’s input-output
 126 mapping after t epochs of training is $\mathbf{h}(t)\mathbf{W}(t)$. Assuming that the network weights align to the
 127 singular vectors of the dataset from early in training, as described by the “silent alignment effect”
 128 [60], we perform a change of variables and write the network mapping in terms of the dataset singular
 129 vectors:

$$\mathbf{h}(t)\mathbf{W}(t) = \omega(t)\mathbf{v}^T, \quad (2)$$

130 where $\omega(t)$ is the network pathway’s scalar effective singular value and the only time-dependent
 131 component of the decomposed network mapping. While the alignment assumption is strong, linear
 132 paradigms with these assumptions have been used successfully in the past [47, 61, 62, 48, 35]. With
 133 the change of variables we can now obtain a closed form equation describing how ω evolves through
 134 time as:

$$\omega(t) = \frac{\lambda}{2} \sinh \left\{ 2 \tanh^{-1} \left[\frac{K \left(C \exp \left(\frac{\text{sgn}(\lambda)K}{\tau} t \right) - 1 \right) - \lambda D \left(C \exp \left(\frac{\text{sgn}(\lambda)K}{\tau} t \right) + 1 \right)}{2S \left(C \exp \left(\frac{\text{sgn}(\lambda)K}{\tau} t \right) + 1 \right)} \right] \right\} \quad (3)$$

135 where C is a defined constant, $\tau = \frac{1}{\eta}$ is the learning time constant and $K = \sqrt{4S^2 + \lambda^2 D^2}$. Eq. 3
 136 shows that K is the variable interacting with time (t) and as a consequence determines how quickly
 137 the network will learn. Three factors affect K fastening learning: 1. S the input-output correlation
 138 matrix singular value, 2. D the input correlation matrix singular value, and 3. $\lambda = h^2 - W^2$ which
 139 denotes the imbalance between the weights of the network. Notice that—as shown in Appendix A— λ is
 140 a conserved quantity and constant throughout training. Thus, given a dataset—which characterises the
 141 S and D matrices—the only property which can promote faster learning in the network is to increase
 142 the imbalance parameter. For our experiments we whiten the input data \mathbf{x} such that $K = \sqrt{4S^2 + \lambda^2}$
 143 to remove one of the interactions within K .

144 Fig. 3(a-c) show a confirmation of the validity our theory by comparing with simulations. Instead,
 145 Fig. 3d represents the main result of this section. We consider both network pathways and vary
 146 the weight imbalance for each (λ_{low} for the pathway with the lower imbalance and λ_{high} for the
 147 pathway with the larger imbalance). We place these two values on the axes and in colour depict
 148 in log scale how close the slower pathway comes to reaching its escaping time across its training.
 149 When negative, it means that during training there is a timestep where the network is less than one

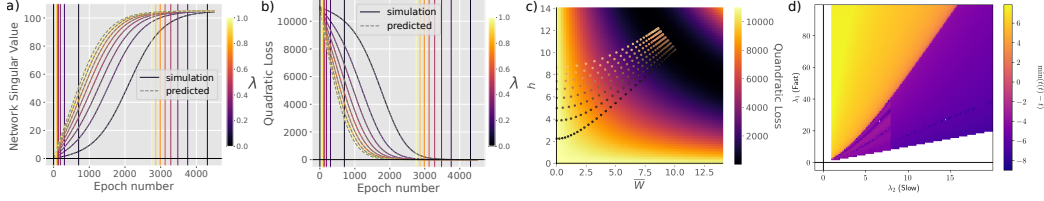


Figure 3: **Linear Dynamics from imbalanced initialisation leads to specialisation.** *Panels a-c)* Show agreement between our theoretical curves and simulations for the training dynamics of: (a) the network’s singular values, (b) the network’s loss, (c) and the network’s movement in weight space. In (a,b) the colour indicates the singular value used for the input weights, while in (c) the colour represents the loss. *Panel d)* shows a phase diagram representing how pathways with different initial weight imbalances lead to specialisation. The two axis represent the initial singular values associated to the different pathways. The colour represents the amount of time it takes the slower pathways to learn in logscale. We see that the more imbalanced the fast pathway relative to the slower pathway, the more likely the network will specialise. The white region represents when the imbalance is reversed.

150 epoch from its escaping time (so it will learn). In this case there will not be specialisation as both
 151 pathways will learn some part of the input-output mapping. When the colour is positive it means
 152 there will be specialisation as the slower pathway is always at least a full epoch away from learning.
 153 It is important to note that the slower pathway’s escaping time is moving constantly as the faster
 154 pathway accounts for variance in the data. This decreases the input-output singular value in K for
 155 this pathway and makes learning slower. Due to this coupling we are also unable to obtain completely
 156 closed form equations for the slower pathway in term’s of the faster pathway’s effective singular
 157 value. However, this phase diagram would not be computationally feasible without the closed-form
 158 escaping time, hitting time and training dynamics (see Appendix C for our process on constructing
 159 this plot). Finally, we only consider imbalances where the output layer is larger than the input layer.
 160 Recent work [33, 35] has shown that having larger input weight pushes the network towards lazy
 161 learning [30] while output heavy imbalance promotes feature learning. Since we are concerned with
 162 the latter in this work, we focus on the output heavy imbalanced setting for both pathways. From
 163 Fig. 3 we see that there is a clear phase transition from non-specialised representations to specialised
 164 ones. This occurs with increasing imbalance of the faster pathway. Increasing the imbalance of the
 165 slower pathway can similarly combat this specialisation pressure. Thus, the relative imbalance of the
 166 two pathway at initialisation will dictate whether specialised representations are learned.

167 3.2 Specialisation underlies disentanglement

168 We extend the results on imbalanced initialisation and applied them, beyond the limited setting
 169 of our framework, in the context of disentangled representation learning, where the goal is to
 170 separate latent factors. [14] introduced the importance of disentanglement for interpretability and
 171 generalisation. A seminal contribution to this domains came with the β -VAE model, where [63]
 172 demonstrated how increasing the KL-divergence term can enforce disentanglement by encouraging
 173 specialised latent representations. Many studies have built upon these foundational frameworks
 174 to enhance disentanglement performance, exploring different training regimes [64, 65] and loss
 175 functions [66, 67, 68]. Here we contribute to this literature by applying our theoretical insights and
 176 examining the impact of initialisation on disentanglement performance.

177 Specifically, we examine how initialisation impacts specialisation in disentanglement learning on
 178 the 3DShapes dataset [69] using the β -VAE model—widely adopted for such tasks [63, 70]. We
 179 implement a β -VAE model, employing the "DeepGaussianLinear" architecture for the decoder
 180 and the "DeepLinear" architecture for the encoder, as specified in [24]. Both architectures are
 181 composed of five fully connected layers with ReLU activations. The model is trained using the Adam
 182 optimiser, optimising a loss function that combines KL divergence and binary cross-entropy-based
 183 reconstruction loss. Additional details are given in Appendix D. In these experiments, we adjust the
 184 variance of the weights in a deep fully-connected encoder, by varying the constant gain of the Xavier
 185 initialisation [71]. Specifically, the first block of layers was initialised with gain g while the readout
 186 layer received a gain $1/g$. Notice that $g = 1$ represents the standard initialisation scheme.

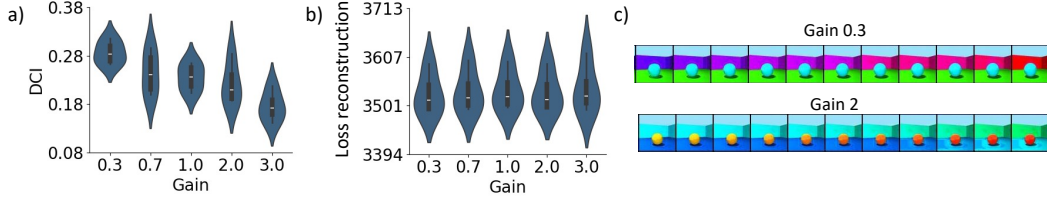


Figure 4: *Panel a)* Violin plots of the DCI values against the gain. *Panel b)* Violin plots of the reconstruction loss against the gain. The standard deviation was computed over four seeds. *Panel c)* Example Traverse of model with gain 2 and .3 respectively showcasing

187 Results are shown in Fig. 4, despite very similar levels of reconstruction loss, networks initialised
 188 with smaller gains improved disentanglement in the β -VAE network, as reflected in higher Disentan-
 189 glement, Completeness, and Informativeness (DCI) scores [72]. This result confirms that modulating
 190 the initialisation gain can either enhance or reduce the network’s disentanglement. Although the scope
 191 of these experiments is limited, they provide preliminary validation of our theoretical framework in
 192 more realistic contexts, encouraging further investigation into alternative initialisation schemes with
 193 varying levels of balance.

194 4 Continual Learning

195 As [15] noted, multi-task learning benefits significantly from task-specific specialisation, allowing the
 196 network to better preserve performance across multiple domains. In the context of continual learning,
 197 [41] and [42] observed that forgetting does not monotonically increase with task similarity. [43]
 198 provided a mechanistic explanation, showing that this phenomenon is due to the interplay between
 199 re-use of specialised neurons and activation of unused ones. In this section, we build on these findings
 200 and show that this phenomenology can be disrupted by initialisation schemes that disincentives
 201 specialisation.

202 4.1 Continual Learning in the two-layer teacher-student setup

203 We use a teacher-student framework, introduced in Sec. 2, which has been analysed in [42, 43]. This
 204 model consists of two randomly initialised teacher networks—one for an upstream task and one for a
 205 downstream task. Each teacher is represented by two-layer neural networks with P^* hidden units and
 206 weights $\mathbf{W}_T^{(1)}, \mathbf{h}_T^{(1)}$ for the upstream task, and $\mathbf{W}_T^{(2)}, \mathbf{h}_T^{(2)}$ for the downstream task. Given a random
 207 input $\mathbf{x} \in \mathbb{R}^d$, drawn i.i.d. from a Gaussian distribution $x_i \sim \mathcal{N}(0, 1)$, the teachers generate labels
 208 according to the equation:

$$y^{(t)} = \mathbf{h}_T^{(t)} \cdot \phi \left(\frac{\mathbf{W}_T^{(t)} \mathbf{x}}{\sqrt{d}} \right) \quad \text{for } t = 1, 2, \quad (4)$$

209 where ϕ is a non-linear activation function, chosen here as $\phi(z) = \text{erf}(z/\sqrt{2})$. This setup allows us
 210 to generate two datasets $\mathcal{D}^{(1)}$ and $\mathcal{D}^{(2)}$, with controlled similarity between the tasks by manipulating
 211 the teacher weights. Specifically, we generate $\mathbf{W}_T^{(1)}, \mathbf{h}_T^{(1)}$, and $\mathbf{h}_T^{(2)}$ with i.i.d. Gaussian entries, while
 212 $\mathbf{W}_T^{(2)}$ is generated as:

$$\mathbf{W}_T^{(2)} = \gamma \mathbf{W}_T^{(1)} + \sqrt{1 - \gamma^2} \mathbf{W}_T^{(\text{aux})}, \quad (5)$$

213 where $\mathbf{W}_T^{(\text{aux})}$ is an auxiliary weight matrix, and γ controls the correlation between tasks. The student
 214 is a two-layer neural network with P hidden units, using the same non-linearity ϕ . It is trained using
 215 online stochastic gradient descent on a squared error loss, with a shared first-layer weight matrix
 216 \mathbf{W} and task-specific readout weights $\mathbf{h}^{(1)}$ and $\mathbf{h}^{(2)}$. For both layers, the initial weights are sampled
 217 i.i.d. from a Gaussian distribution, with the first-layer weights \mathbf{W} having standard deviation σ_W .
 218 While most previous studies follow a similar scheme for the readout weights, we introduce a novel
 219 initialisation scheme using polar coordinates, as detailed in Eq. 11. The updates for \mathbf{W} and $\mathbf{h}^{(t)}$ at

220 iteration e , under SGD on the squared error loss, are given by:

$$\mathbf{W}[e+1] = \mathbf{W}[e] - \frac{\eta}{\sqrt{d}} \left(\mathbf{h}^{(t)} \cdot \phi \left(\frac{\mathbf{W}\mathbf{x}}{\sqrt{d}} \right) - y^{(t)} \right) \phi' \left(\frac{\mathbf{W}\mathbf{x}}{\sqrt{d}} \right) \mathbf{v}^{(t)} \mathbf{x}, \quad (6)$$

$$\mathbf{h}^{(t)}[e+1] = \mathbf{h}^{(t)}[e] - \frac{\eta}{d} \left(\mathbf{h}^{(t)} \cdot \phi \left(\frac{\mathbf{W}\mathbf{x}}{\sqrt{d}} \right) - y^{(t)} \right) \phi \left(\frac{\mathbf{W}\mathbf{x}}{\sqrt{d}} \right), \quad (7)$$

221 where η is the learning rate and $y^{(t)}$ is the target output from the teacher network for task t .

222 In the large input dimension limit $d \rightarrow \infty$, key observables, such as the generalisation error, can be
223 captured by a few order parameters:

$$\mathbf{Q} = \frac{1}{d} \mathbf{W}\mathbf{W}^T, \quad \mathbf{R}^{(t)} = \frac{1}{d} \mathbf{W}\mathbf{W}_T^{(t),T}, \quad \mathbf{T}^{(t,t')} = \frac{1}{d} \mathbf{W}_T^{(t)} \mathbf{W}_T^{(t'),T}, \quad \mathbf{h}^{(t)}, \quad \mathbf{h}_T^{(t)}; \quad (8)$$

224 where $t, t' \in \{1, 2\}$ refer to the two tasks. The generalisation error for task t is then:

$$\begin{aligned} \epsilon^{(t)} &= \frac{1}{2} \mathbb{E}_{\mathbf{x}} \left[\left(\mathbf{h}^{(t)} \cdot \phi \left(\frac{\mathbf{W}\mathbf{x}}{\sqrt{d}} \right) - y^{(t)} \right)^2 \right] \\ &= I_{21}(\mathbf{Q}, \mathbf{h}^{(t)}) + I_{21}(\mathbf{T}^{(t,t)}, \mathbf{h}_T^{(t)}) - \frac{1}{2} I_{22}(\mathbf{Q}, \mathbf{R}^{(t)}, \mathbf{T}^{(t,t)}, \mathbf{h}^{(t)}, \mathbf{h}_T^{(t)}), \end{aligned} \quad (9)$$

225 where I_{21} and I_{22} are explicit functions of the order parameters, detailed in Appendix B. The
226 evolution of these parameters throughout training can be tracked to study the learning dynamics, as
227 first shown in [39, 40, 45]. For the specific case of continual learning, [42] derived the governing
228 ordinary differential equations (ODEs), provided in Appendix B.

229 4.2 Specialisation relevance for continual learning

230 The continual learning results in the teacher-student setup, including the non-monotonic relationship
231 between catastrophic forgetting and task similarity, often implicitly assume that the student has
232 specialised to the teacher in the first task. This assumption allows for spare capacity to represent the
233 second task. However, as shown in Fig. 1b, there are regimes where this assumption of specialisation
234 is violated. Here, we expand on these findings and their implications for forgetting.

235 A student can effectively ignore a unit in two ways: either the unit's post-activation is near 0 (inactive),
236 or the corresponding second-layer weight is 0. This motivates three measures for specialisation based
237 on the definition of entropy—over the hidden units, head weights, and the product of both:

$$H_h = - \sum_i^P \tilde{h}_i \log |\tilde{h}_i|, \quad H_Q = - \sum_i^P \tilde{Q}_{ii} \log \tilde{Q}_{ii}, \quad H_m = - \sum_i^P \tilde{Q}_{ii} |\tilde{h}_i| \log(\tilde{Q}_{ii} |\tilde{h}_i|); \quad (10)$$

238 where the tilde denote normalisation, i.e. $|\tilde{h}_i| = \frac{|h_i|}{\sum_i^P |h_i|}$ and $\tilde{Q}_{ii} = \frac{Q_{ii}}{\sum_i^P Q_{ii}}$. Maximum entropy in
239 these measures corresponds to no specialisation, while minimum entropy corresponds to maximum
240 specialisation.

241 We can investigate how these measures vary as a function of different properties of the problem setup,
242 in particular those related to initialisation. To simplify the analysis, we begin with the case where
243 the optimal number of tasks is $P^* = 1$ and the network has $P = 2$ output units. This allows us to
244 initialise the second layer weights in polar coordinates, with precise and interpretable control over
245 scale and asymmetry of weights. Formally we parameterise our readout initialisations according to

$$\mathbf{h}^{(t)}[0; r^{(t)}, \theta^{(t)}] = (r^{(t)} \cos \theta^{(t)}, r^{(t)} \sin \theta^{(t)}). \quad (11)$$

246 Fig. 5 contain phase diagrams showing how the entropy measures in Eq. 10 vary with the initialisation
247 parameters $r^{(t)}$, $\theta^{(t)}$, and σ_W . We can make several observations: (i) the strongest determinant of
248 specialisation is the asymmetry in the second layer weights, i.e. the θ parameter. (ii) this is the
249 case for both ReLU and sigmoidal activation functions, reinforcing the point made in the example
250 from Fig. 1b. (iii) the scale of initialisations (parameters σ_W , r) are also important.

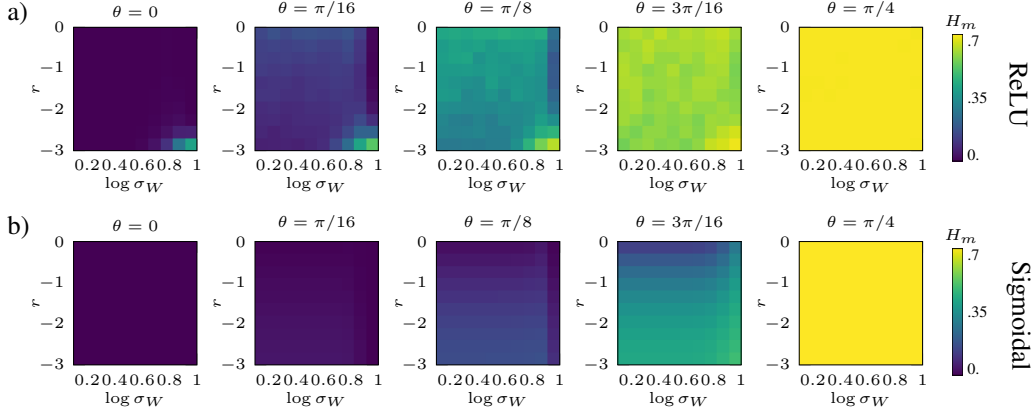


Figure 5: **Phase diagrams show significance of initialisation for specialisation.** The phase diagrams show with colour the aggregated entropy Eq. 10 evaluated for different initialisations. On the x-axis we span over the standard deviation of the first layer. The second layer is initialised using polar coordinates, and the y-axis represents the norm while the different panels give the angle spanning from orthogonal units ($\theta = 0$) to identical units ($\theta = \pi/4$). Specialisation is achieved by blue-leaning initialisations, while yellow-leaning ones exhibit high entropy and therefore non-specialised solutions. Additional results can be found in Appendix E. standard

251 4.3 Specialisation underlies Maslow’s hammer

252 The phase diagrams in Fig. 5 demonstrate that initialisation can drastically change the type of
 253 solutions found by the student after training on one teacher. While this may be inconsequential if the
 254 generalisation error remains unaffected, in many cases, the precise nature of the learned representation
 255 can significantly impact downstream tasks.

256 In the worst case scenario, the student undergoes no specialisation during the first task. During the
 257 second task there is no notion of the trade-off between node re-use and node activation discussed
 258 in [43]; rather the student continues to find a non-specialised solution to the second teacher, effectively
 259 fully re-using it’s entire representation for the second task. Consequently, the amount of forgetting
 260 with respect to the initial task decreases monotonically with task similarity, thereby breaking the
 261 U-shaped pattern characteristic of Maslow’s hammer that has been observed in various continual
 262 learning setups [41]. This extreme case is illustrated in Fig. 6. Further, *even with* specialisation
 263 after the first task, large asymmetric initialisation in the second task readout weights can induce this
 264 monotonic relationship, again by pushing the student into re-use rather than activation.

265 In a broader context, a rich diversity of behaviours can emerge, driven by factors such as the
 266 initialisation schemes, the scale of weights in the first layer, and the readout heads for both tasks.
 267 A glimpse of this behavioural diversity is provided in Appendix F, where we further explore the
 268 interaction between these factors and their impact on forgetting in continual learning.

269 4.4 Specialisation underlies EWC

270 The findings relating specialisation to forgetting from subsection 4.3 have direct consequences for
 271 interference mitigation strategies such as EWC. EWC is a regularisation-based method that computes
 272 a measure of ‘importance’ for each weight with respect to a task via the Fischer information [21].
 273 Subsequently a squared penalty scaled by this importance is applied to deviation of this weight during
 274 learning of future tasks as follows:

$$\mathcal{L}_{\text{EWC}}(\mathbf{W}) = \mathcal{L}(\mathbf{W}) + \frac{\xi}{2} \sum_i F_i (W_i - W_i^*), \quad (12)$$

275 where F is the Fischer information matrix, ξ is a regularisation strength parameter, and \mathbf{W}^* are the
 276 weights at the end of training on the first task.

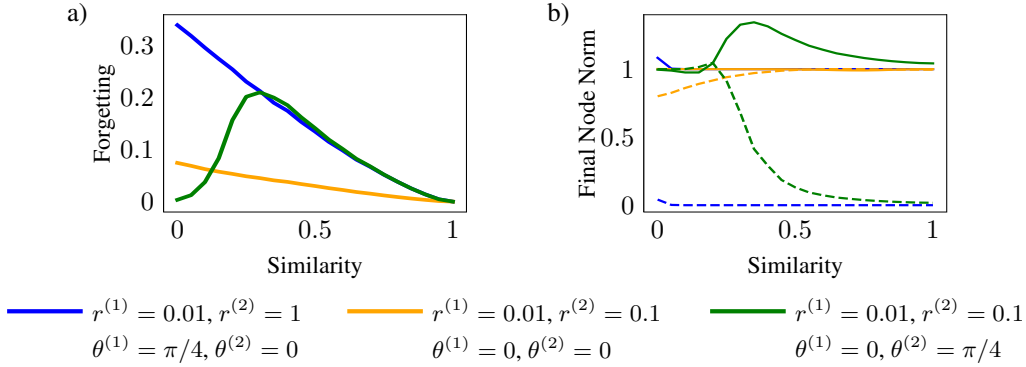


Figure 6: **Initialisation and specialisation properties can influence profile of forgetting vs. similarity.** (a) forgetting as a function of task similarity can be both monotonic, shown here for the cases of specialisation after the first task + large second head initialisation (blue), and no-specialisation during both tasks (orange); or non-monotonic (green, as characterised by Maslow’s hammer [43]). (b) the final norm of the two nodes (one solid and one dashed), i.e. at the end of training on both tasks, as a function of task similarity. In the cases that lead to monotonic forgetting, nodes are fully re-used, either because the corresponding new head is initialised large (orange) or because the new head is symmetrically initialised and the nodes continue to represent redundant information during the second task (blue). *Params:* $N = 10000, \eta = 1, K^* = 1, K = 2, \sigma_w = 0.001$.

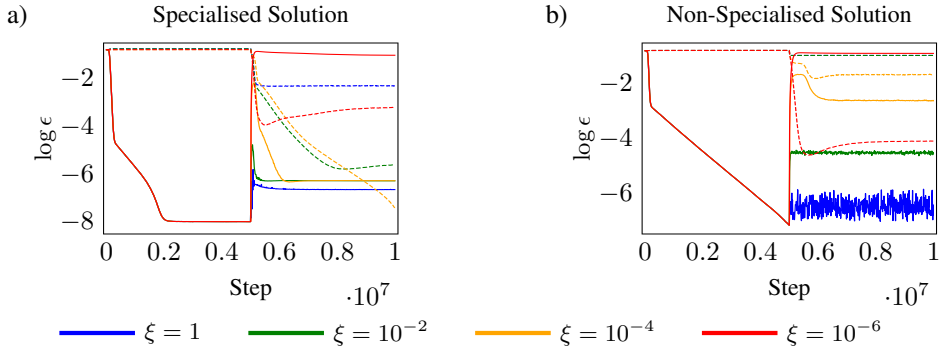


Figure 7: **EWC is strongly reliant on specialisation.** We show the generalisation error in the first (solid line) and second (dashed) task for different EWC regularisation strengths. (a) When the student finds a specialised solution to the first task, there is a range of EWC regularisation strength ξ for which the activated units can remain fixed and spare capacity can be used to learn the second task—leading to low generalisation error in both tasks ($\xi = 10^{-2}, \xi = 10^{-4}$ perform very well). (b) When the student does not specialise in the first task, EWC reduces to an inflexible regulariser that either penalises plasticity everywhere—leading to little forgetting but no further learning (e.g. $\xi = 1$), or does not penalise any plasticity—leading to catastrophic forgetting (e.g. $\xi = 10^{-6}$).

277 In cases where the network does not specialise, i.e. multiple student nodes learn redundant repre-
 278 sentations for a given teacher node, the nodes have equal importance. Consequently EWC cannot
 279 distinguish between these sets of weights and depending on the regularisation parameter λ either lets
 280 these nodes move during training on the second task (under-regularises) leading to forgetting, or lets
 281 none move (over-regularises) leading to no transfer. We show results illustrating this behaviour in
 282 the teacher-student setup in Fig. 7. In particular we show the regime of intermediate task similarity,
 283 wherein [43] previously argued that EWC should perform better than methods such as replay.

284 5 Limitations and Perspectives

285 This work operates within simplified frameworks, which—while widely used in the analysis of neural
286 networks—do not fully capture the complexity of modern architectures and real-world data. Our
287 experiments rely on Gaussian input data and simplified input-output relations, which are far removed
288 from the intricacies of real-world scenarios. A natural next step is to extend our analysis to more
289 realistic generative models, such as the hidden manifold model [73] or the superstatistical generative
290 model [74], which offer more structured data distributions and better capture observations from real
291 data experiments.

292 Another promising direction is to complement analytical approaches with numerical experiments on
293 controlled real-world datasets. While this may sacrifice some analytical tractability, it brings us closer
294 to addressing practical challenges. For instance, transfer learning settings, such as those explored in
295 [75], provide a useful benchmark for testing our theoretical findings in more complex environments.

296 While the current work remains theoretical in nature, focusing on simplified models for analytical
297 tractability, a thorough exploration of the practical implications of our findings, particularly in
298 disentangled representation learning, is beyond the scope of this paper. However, we aim to address
299 this in future work by shifting towards a more experimental approach. Specifically, we plan to
300 explore a broader range of network architectures, datasets—such as Car3D [?] and dSprites [76]—and
301 evaluation metrics—such as SAP [68, 63]. This future study will allow us to validate our theoretical
302 insights and fully assess their relevance in real-world settings.

303 References

- 304 [1] Horace B Barlow. Single units and sensation: a neuron doctrine for perceptual psychology?
305 *Perception*, 1(4):371–394, 1972.
- 306 [2] John J Hopfield. Neural networks and physical systems with emergent collective computational
307 abilities. *Proceedings of the national academy of sciences*, 79(8):2554–2558, 1982.
- 308 [3] Colin Blakemore, James PJ Muncey, and Rosalind M Ridley. Stimulus specificity in the human
309 visual system. *Vision research*, 13(10):1915–1931, 1973.
- 310 [4] R Quian Quiroga, Leila Reddy, Gabriel Kreiman, Christof Koch, and Itzhak Fried. Invariant
311 visual representation by single neurons in the human brain. *Nature*, 435(7045):1102–1107,
312 2005.
- 313 [5] Apostolos P Georgopoulos, Andrew B Schwartz, and Ronald E Kettner. Neuronal population
314 coding of movement direction. *Science*, 233(4771):1416–1419, 1986.
- 315 [6] Alumi Ishai, Leslie G Ungerleider, Alex Martin, and James V Haxby. The representation
316 of objects in the human occipital and temporal cortex. *Journal of cognitive neuroscience*,
317 12(Supplement 2):35–51, 2000.
- 318 [7] Bruno B Averbeck, Peter E Latham, and Alexandre Pouget. Neural correlations, population
319 coding and computation. *Nature reviews neuroscience*, 7(5):358–366, 2006.
- 320 [8] Yann LeCun, Bernhard Boser, John S Denker, Donnie Henderson, Richard E Howard, Wayne
321 Hubbard, and Lawrence D Jackel. Backpropagation applied to handwritten zip code recognition.
322 *Neural computation*, 1(4):541–551, 1989.
- 323 [9] Dumitru Erhan, Aaron Courville, Yoshua Bengio, and Pascal Vincent. Why does unsupervised
324 pre-training help deep learning? In *Proceedings of the thirteenth international conference
325 on artificial intelligence and statistics*, pages 201–208. JMLR Workshop and Conference
326 Proceedings, 2010.
- 327 [10] Jason Yosinski, Jeff Clune, Yoshua Bengio, and Hod Lipson. How transferable are features in
328 deep neural networks? *Advances in neural information processing systems*, 27, 2014.
- 329 [11] Matthew D Zeiler and Rob Fergus. Visualizing and understanding convolutional networks. In
330 *Computer Vision—ECCV 2014: 13th European Conference, Zurich, Switzerland, September
331 6–12, 2014, Proceedings, Part I 13*, pages 818–833. Springer, 2014.

- 332 [12] Elena Voita, David Talbot, Fedor Moiseev, Rico Sennrich, and Ivan Titov. Analyzing multi-
333 head self-attention: Specialized heads do the heavy lifting, the rest can be pruned. In Anna
334 Korhonen, David R. Traum, and Lluís Màrquez, editors, *Proceedings of the 57th Conference of*
335 *the Association for Computational Linguistics, ACL 2019, Florence, Italy, July 28- August 2,*
336 *2019, Volume 1: Long Papers*, pages 5797–5808. Association for Computational Linguistics,
337 2019.
- 338 [13] Adly Templeton. *Scaling monosemanticity: Extracting interpretable features from claude 3*
339 *sonnet*. Anthropic, 2024.
- 340 [14] Yoshua Bengio, Aaron Courville, and Pascal Vincent. Representation learning: A review and
341 new perspectives. *IEEE transactions on pattern analysis and machine intelligence*, 35(8):1798–
342 1828, 2013.
- 343 [15] Rich Caruana. Multitask learning. *Machine learning*, 28:41–75, 1997.
- 344 [16] Michael McCloskey and Neal J Cohen. Catastrophic interference in connectionist networks:
345 The sequential learning problem. In *Psychology of learning and motivation*, volume 24, pages
346 109–165. Elsevier, 1989.
- 347 [17] Roger Ratcliff. Connectionist models of recognition memory: constraints imposed by learning
348 and forgetting functions. *Psychological review*, 97(2):285, 1990.
- 349 [18] Steven C Sudderth and YL Kergosien. Rule-injection hints as a means of improving network
350 performance and learning time. In *European association for signal processing workshop*, pages
351 120–129. Springer, 1990.
- 352 [19] German I Parisi, Ronald Kemker, Jose L Part, Christopher Kanan, and Stefan Wermter. Continual
353 lifelong learning with neural networks: A review. *Neural networks*, 113:54–71, 2019.
- 354 [20] Matthias De Lange, Rahaf Aljundi, Marc Masana, Sarah Parisot, Xu Jia, Aleš Leonardis,
355 Gregory Slabaugh, and Tinne Tuytelaars. A continual learning survey: Defying forgetting in
356 classification tasks. *IEEE transactions on pattern analysis and machine intelligence*, 44(7):3366–
357 3385, 2021.
- 358 [21] James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins,
359 Andrei A Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al.
360 Overcoming catastrophic forgetting in neural networks. *Proceedings of the national academy of*
361 *sciences*, 114(13):3521–3526, 2017.
- 362 [22] Friedemann Zenke, Ben Poole, and Surya Ganguli. Continual learning through synaptic
363 intelligence. In *International conference on machine learning*, pages 3987–3995. PMLR, 2017.
- 364 [23] Zhizhong Li and Derek Hoiem. Learning without forgetting. *IEEE transactions on pattern*
365 *analysis and machine intelligence*, 40(12):2935–2947, 2017.
- 366 [24] Francesco Locatello, Stefan Bauer, Mario Lucic, Gunnar Raetsch, Sylvain Gelly, Bernhard
367 Schölkopf, and Olivier Bachem. Challenging common assumptions in the unsupervised learning
368 of disentangled representations. In *international conference on machine learning*, pages 4114–
369 4124. PMLR, 2019.
- 370 [25] Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and
371 generalization in neural networks. *Advances in neural information processing systems*, 31,
372 2018.
- 373 [26] Song Mei, Andrea Montanari, and Phan-Minh Nguyen. A mean field view of the landscape of
374 two-layer neural networks. *Proceedings of the National Academy of Sciences*, 115(33):E7665–
375 E7671, 2018.
- 376 [27] Lenaïc Chizat and Francis Bach. On the global convergence of gradient descent for over-
377 parameterized models using optimal transport. *Advances in neural information processing*
378 *systems*, 31, 2018.

- 379 [28] Grant M Rotskoff and Eric Vanden-Eijnden. Neural networks as interacting particle systems:
380 Asymptotic convexity of the loss landscape and universal scaling of the approximation error.
381 *stat*, 1050:22, 2018.
- 382 [29] Lenaic Chizat, Edouard Oyallon, and Francis Bach. On lazy training in differentiable program-
383 ming. *Advances in neural information processing systems*, 32, 2019.
- 384 [30] Mario Geiger, Stefano Spigler, Arthur Jacot, and Matthieu Wyart. Disentangling feature and
385 lazy training in deep neural networks. *Journal of Statistical Mechanics: Theory and Experiment*,
386 2020(11):113301, 2020.
- 387 [31] Blake Bordelon and Cengiz Pehlevan. Self-consistent dynamical field theory of kernel evolution
388 in wide neural networks. *Advances in Neural Information Processing Systems*, 35:32240–32256,
389 2022.
- 390 [32] Salma Tarmoun, Guilherme Franca, Benjamin D Haeffele, and Rene Vidal. Understanding the
391 dynamics of gradient flow in overparameterized linear models. In *International Conference on*
392 *Machine Learning*, pages 10153–10161. PMLR, 2021.
- 393 [33] Daniel Kunin, Allan Raventós, Clémentine Dominé, Feng Chen, David Klindt, Andrew Saxe,
394 and Surya Ganguli. Get rich quick: exact solutions reveal how unbalanced initializations
395 promote rapid feature learning, 06 2024.
- 396 [34] Yizhou Xu and Liu Ziyin. When does feature learning happen? perspective from an analytically
397 solvable model. *arXiv preprint arXiv:2401.07085*, 2024.
- 398 [35] Clémentine C. J. Dominé, Nicolas Anguita, Alexandra M. Proca, Lukas Braun, Daniel Kunin,
399 Pedro A. M. Mediano, and Andrew M. Saxe. From lazy to rich: Exact learning dynamics in
400 deep linear networks, 2024.
- 401 [36] Aditya Vardhan Varre, Maria-Luiza Vladarean, Loucas Pillaud-Vivien, and Nicolas Flammarion.
402 On the spectral bias of two-layer linear networks. *Advances in Neural Information Processing*
403 *Systems*, 36, 2024.
- 404 [37] Andrew M Saxe, James L McClelland, and Surya Ganguli. Exact solutions to the nonlinear
405 dynamics of learning in deep linear neural networks. *arXiv preprint arXiv:1312.6120*, 2013.
- 406 [38] David Saad and Sara A Solla. On-line learning in soft committee machines. *Physical Review E*,
407 52(4):4225, 1995.
- 408 [39] David Saad and Sara Solla. Dynamics of on-line gradient descent learning for multilayer neural
409 networks. *Advances in neural information processing systems*, 8, 1995.
- 410 [40] Michael Biehl and Holm Schwarze. Learning by on-line gradient descent. *Journal of Physics A:*
411 *Mathematical and general*, 28(3):643, 1995.
- 412 [41] Vinay V Ramasesh, Ethan Dyer, and Maithra Raghu. Anatomy of catastrophic forgetting:
413 Hidden representations and task semantics. *arXiv preprint arXiv:2007.07400*, 2020.
- 414 [42] Sebastian Lee, Sebastian Goldt, and Andrew Saxe. Continual learning in the teacher-student
415 setup: Impact of task similarity. In *International Conference on Machine Learning*, pages
416 6109–6119. PMLR, 2021.
- 417 [43] Sebastian Lee, Stefano Sarao Mannelli, Claudia Clopath, Sebastian Goldt, and Andrew Saxe.
418 Maslow’s hammer for catastrophic forgetting: Node re-use vs node activation. *arXiv preprint*
419 *arXiv:2205.09029*, 2022.
- 420 [44] Elizabeth Gardner and Bernard Derrida. Three unfinished works on the optimal storage capacity
421 of networks. *Journal of Physics A: Mathematical and General*, 22(12):1983, 1989.
- 422 [45] Sebastian Goldt, Madhu Advani, Andrew M Saxe, Florent Krzakala, and Lenka Zdeborová.
423 Dynamics of stochastic gradient descent for two-layer neural networks in the teacher-student
424 setup. *Advances in neural information processing systems*, 32, 2019.

- 425 [46] Andrew Saxe, Shagun Sodhani, and Sam Jay Lewallen. The neural race reduction: Dynamics
426 of abstraction in gated networks. In *International Conference on Machine Learning*, pages
427 19287–19309. PMLR, 2022.
- 428 [47] Andrew M Saxe, James L McClelland, and Surya Ganguli. A mathematical theory of semantic
429 development in deep neural networks. *Proceedings of the National Academy of Sciences*,
430 116(23):11537–11546, 2019.
- 431 [48] Devon Jarvis, Richard Klein, Benjamin Rosman, and Andrew M Saxe. On the specialization of
432 neural modules. In *The Eleventh International Conference on Learning Representations*, 2023.
- 433 [49] Jacob Andreas, Marcus Rohrbach, Trevor Darrell, and Dan Klein. Neural module networks. In
434 *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 39–48,
435 2016.
- 436 [50] Ronghang Hu, Jacob Andreas, Marcus Rohrbach, Trevor Darrell, and Kate Saenko. Learning
437 to reason: End-to-end module networks for visual question answering. In *Proceedings of the
438 IEEE International Conference on Computer Vision*, pages 804–813, 2017.
- 439 [51] Ronghang Hu, Jacob Andreas, Trevor Darrell, and Kate Saenko. Explainable neural computation
440 via stack neural module networks. In *Proceedings of the European conference on computer
441 vision (ECCV)*, pages 53–69, 2018.
- 442 [52] Jacob Andreas. Measuring compositionality in representation learning. In *International
443 Conference on Learning Representations*, 2018.
- 444 [53] Saeed Masoudnia and Reza Ebrahimpour. Mixture of experts: a literature survey. *Artificial
445 Intelligence Review*, 42(2):275–293, 2014.
- 446 [54] Emmanuel Bengio, Pierre-Luc Bacon, Joelle Pineau, and Doina Precup. Conditional computa-
447 tion in neural networks for faster models. *arXiv preprint arXiv:1511.06297*, 2015.
- 448 [55] Noam Shazeer, Azalia Mirhoseini, Krzysztof Maziarz, Andy Davis, Quoc Le, Geoffrey Hinton,
449 and Jeff Dean. Outrageously large neural networks: The sparsely-gated mixture-of-experts
450 layer. *arXiv preprint arXiv:1701.06538*, 2017.
- 451 [56] Paul Smolensky, R Thomas McCoy, Roland Fernandez, Matthew Goldrick, and Jianfeng Gao.
452 Neurocompositional computing: From the central paradox of cognition to a new generation of
453 ai systems. *arXiv preprint arXiv:2205.01128*, 2022.
- 454 [57] Michael B Chang, Abhishek Gupta, Sergey Levine, and Thomas L Griffiths. Automati-
455 cally composing representation transformations as a means for generalization. *arXiv preprint
456 arXiv:1807.04640*, 2018.
- 457 [58] Anirudh Goyal, Alex Lamb, Jordan Hoffmann, Shagun Sodhani, Sergey Levine, Yoshua Bengio,
458 and Bernhard Schölkopf. Recurrent independent mechanisms. *arXiv preprint arXiv:1909.10893*,
459 2019.
- 460 [59] Ishita Dasgupta, Erin Grant, and Tom Griffiths. Distinguishing rule and exemplar-based
461 generalization in learning systems. In *International Conference on Machine Learning*, pages
462 4816–4830. PMLR, 2022.
- 463 [60] Alexander Atanasov, Blake Bordelon, and Cengiz Pehlevan. Neural networks as kernel learners:
464 The silent alignment effect. *arXiv preprint arXiv:2111.00034*, 2021.
- 465 [61] A.K. Lampinen and S. Ganguli. An analytic theory of generalization dynamics and transfer
466 learning in deep linear networks. In T. Sainath, editor, *International Conference on Learning
467 Representations*, 2019. arXiv: 1809.10374.
- 468 [62] Lukas Braun, Clémentine Dominé, James Fitzgerald, and Andrew Saxe. Exact learning dynam-
469 ics of deep linear networks with prior knowledge. *Advances in Neural Information Processing
470 Systems*, 35:6615–6629, 2022.

- 471 [63] Irina Higgins, Loic Matthey, Arka Pal, Christopher P Burgess, Xavier Glorot, Matthew M
472 Botvinick, Shakir Mohamed, and Alexander Lerchner. beta-vae: Learning basic visual concepts
473 with a constrained variational framework. *ICLR (Poster)*, 3, 2017.
- 474 [64] Francesco Locatello, Ben Poole, Gunnar Rätsch, Bernhard Schölkopf, Olivier Bachem, and
475 Michael Tschannen. Weakly-supervised disentanglement without compromises, 2020.
- 476 [65] Marco Fumero, Luca Cosmo, Simone Melzi, and Emanuele Rodola. Learning disentangled
477 representations via product manifold projection. In Marina Meila and Tong Zhang, editors, *Pro-
478 ceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings
479 of Machine Learning Research*, pages 3530–3540. PMLR, 18–24 Jul 2021.
- 480 [66] Ricky T. Q. Chen, Xuechen Li, Roger Grosse, and David Duvenaud. Isolating sources of
481 disentanglement in variational autoencoders, 2019.
- 482 [67] Hyunjik Kim and Andriy Mnih. Disentangling by factorising, 2019.
- 483 [68] Abhishek Kumar, Prasanna Sattigeri, and Avinash Balakrishnan. Variational inference of
484 disentangled latent concepts from unlabeled observations, 2018.
- 485 [69] Chris Burgess and Hyunjik Kim. 3d shapes dataset. [https://github.com/deepmind/3dshapes-
486 dataset/](https://github.com/deepmind/3dshapes-dataset/), 2018.
- 487 [70] Christopher P Burgess, Irina Higgins, Arka Pal, Loic Matthey, Nick Watters, Guillaume Des-
488 jardins, and Alexander Lerchner. Understanding disentangling in β -vae. *arXiv preprint
489 arXiv:1804.03599*, 2018.
- 490 [71] Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedfor-
491 ward neural networks. In *Proceedings of the thirteenth international conference on artificial
492 intelligence and statistics*, pages 249–256. JMLR Workshop and Conference Proceedings, 2010.
- 493 [72] Cian Eastwood and Christopher KI Williams. A framework for the quantitative evaluation of
494 disentangled representations. In *6th International Conference on Learning Representations*,
495 2018.
- 496 [73] Sebastian Goldt, Marc Mézard, Florent Krzakala, and Lenka Zdeborová. Modeling the influence
497 of data structure on learning in neural networks: The hidden manifold model. *Physical Review
498 X*, 10(4):041044, 2020.
- 499 [74] Urte Adomaityte, Gabriele Sicuro, and Pierpaolo Vivo. Classification of superstatistical features
500 in high dimensions. In *2023 Conference on Neural Information Procecessing Systems*, 2023.
- 501 [75] Federica Gerace, Diego Doimo, Stefano Sarao Mannelli, Luca Saglietti, and Alessandro Laio.
502 How to choose the right transfer learning protocol? a qualitative analysis in a controlled set-up.
503 *Transactions on Machine Learning Research*, 2024.
- 504 [76] Loic Matthey, Irina Higgins, Demis Hassabis, and Alexander Lerchner. dsprites: Disentangle-
505 ment testing sprites dataset, 2017.
- 506 [77] Anchit Jain, Rozhin Nobahari, Aristide Baratin, and Stefano Sarao Mannelli. Bias in motion:
507 Theoretical insights into the dynamics of bias in sgd training. *arXiv preprint arXiv:2405.18296*,
508 2024.
- 509 [78] Gerard Ben Arous, Reza Gheissari, and Aukosh Jagannath. High-dimensional limit theorems
510 for sgd: Effective dynamics and critical scaling. *Advances in Neural Information Processing
511 Systems*, 35:25349–25362, 2022.
- 512 [79] Amir H. Abdi, Purang Abolmaesumi, and Sidney Fels. Variational learning with
513 disentanglement-pytorch. *arXiv preprint arXiv:1912.05184*, 2019.
- 514 [80] Xiaobiao Du, Haiyang Sun, Shuyun Wang, Zhuojie Wu, Hongwei Sheng, Jiaying Ying, Ming
515 Lu, Tianqing Zhu, Kun Zhan, and Xin Yu. 3drealcar: An in-the-wild rgb-d car dataset with
516 360-degree views, 2024.

517 **A Hyperbolic-Linear Dynamics**

518 Consider a linear network performing a regression task with one hidden layer computing output
 519 $\hat{Y} = hWX$ in response to an input batch of data X , with P datapoints, and trained to minimize the
 520 quadratic loss using gradient descent:

$$L(W, h) = \sum_{i=1}^P \frac{1}{2} \|y_i - hW\mathbf{x}_i\|_2^2$$

521 This gives the learning rules for each layer with learning rate η as:

$$\Delta W = \eta P h^T (\Sigma^{yx} - hW\Sigma^x); \quad \Delta h = \eta P (\Sigma^{yx} - hW\Sigma^x) W^T$$

522 These equations can be derived for a batch of data using the linearity of expectation, where $\Sigma^x =$
 523 $\mathbb{E}[XX^T]$ is the input correlation matrix and $\Sigma^{yx} = \mathbb{E}[YX^T]$ is the input-output correlation matrix,
 524 as follows:

$$\begin{aligned} \Delta W &= \eta \frac{d}{dW} L(W, h) \\ &= \eta \frac{d}{dW} \sum_{i=1}^P \frac{1}{2} (Y_i - hWX_i)^T (Y_i - hWX_i) \\ &= \eta \sum_{i=1}^P h^T (Y_i - hWX_i) X_i^T \\ &= \eta P \frac{1}{P} \sum_{i=1}^P h^T (Y_i - hWX_i) X_i^T \\ &= \eta P \mathbb{E}[h^T (Y_i X_i^T - hWX_i X_i^T)] \\ &= \eta P h^T (\mathbb{E}[Y_i X_i^T] - hW \mathbb{E}[X_i X_i^T]) \\ &= \eta P h^T (\Sigma^{yx} - hW\Sigma^x) \end{aligned}$$

$$\begin{aligned} \Delta h &= \eta \frac{d}{dh} L(W, h) \\ &= \eta \frac{d}{dh} \sum_{i=1}^P \frac{1}{2} (Y_i - hWX_i)^T (Y_i - hWX_i) \\ &= \eta \sum_{i=1}^P (Y_i - hWX_i) (WX_i)^T \\ &= \eta P \frac{1}{P} \sum_{i=1}^P (Y_i - hWX_i) X_i^T W^T \\ &= \eta P \mathbb{E}[(Y_i X_i^T - hWX_i X_i^T)] W^T \\ &= \eta P (\mathbb{E}[Y_i X_i^T] - hW \mathbb{E}[X_i X_i^T]) W^T \\ &= \eta P (\Sigma^{yx} - hW\Sigma^x) W^T \end{aligned}$$

525 By using a small learning rate η and taking the continuous time limit, the mean change in weights is
 526 given by:

$$\tau \frac{d}{dt} W = h^T (\Sigma^{yx} - hW\Sigma^x); \quad \tau \frac{d}{dt} h = (\Sigma^{yx} - hW\Sigma^x) W^T$$

527 where $\tau = \frac{1}{P\eta}$ is the learning time constant. Here, t measures units of learning epochs. It is helpful
 528 to note that since we are using a small learning rate the full batch gradient descent and stochastic
 529 gradient descent dynamics will be the same.

530 [47] has shown that the learning dynamics depend on the singular value decomposition of:

$$\Sigma^{yx} = USV^T = \sum_{\alpha=1}^{r_y} \sigma_{\alpha} u^{\alpha} v^{\alpha T}; \quad \Sigma^x = VDV^T = \sum_{\alpha=1}^{r_x} \delta_{\alpha} u^{\alpha} v^{\alpha T}$$

531 Here r_y and r_x denote the ranks of the matrices. To solve for the dynamics we require that Σ^{yx} and
 532 Σ^x are mutually diagonalizable such that the right singular vectors V of Σ^{yx} are also the singular
 533 vectors of Σ^x . We verify that this is true for the tasks considered in this work and assume it to be
 534 true for these derivations. We also assume that the network has at least r_y hidden neurons (the rank
 535 of Σ^{yx} which determines the number of singular values in the input-output covariance matrix) so
 536 that it can learn the desired mapping perfectly. If this is not the case then the model will learn the
 537 top n_h singular values of the input-output mapping where n_h is the number of hidden neurons [37].
 538 To ease notation for the remainder of this section we will use n_h to denote both the number of
 539 hidden neurons and rank of Σ^{yx} . S and D then are diagonal matrices of the singular values of the
 540 input-output correlation and input correlation matrices respectfully.

541

542 We now perform a change of variables using the SVD of the dataset statistics. The purpose of this
 543 step is to decouple the complex dynamics of the weights of the network, with interacting terms, into
 544 multiple one-dimensional systems. Specifically we set:

$$h = U\bar{h}R^T; \quad W = R\bar{W}V^T$$

545 where R is an arbitrary orthogonal matrix such that $R^T R = I$. Substituting this into the gradient
 546 descent update rules for the parameters above yields:

$$\begin{aligned} \tau \frac{d}{dt} W &= h^T (\Sigma^{yx} - hW\Sigma^x) \\ \tau \frac{d}{dt} (R\bar{W}V^T) &= R\bar{h}U^T (USV^T - U\bar{h}R^T R\bar{W}V^T VDV^T) \\ \tau \frac{d}{dt} (R\bar{W}V^T) &= R\bar{h}(SV^T - \bar{h}\bar{W}DV^T) \\ \tau \frac{d}{dt} \bar{W} &= \bar{h}(S - \bar{h}\bar{W}D) \end{aligned}$$

547 and

$$\begin{aligned} \tau \frac{d}{dt} h &= (\Sigma^{yx} - hW\Sigma^x)W^T \\ \tau \frac{d}{dt} (U\bar{h}R^T) &= (USV^T - U\bar{h}R^T R\bar{W}V^T VDV^T)V\bar{W}R^T \\ \tau \frac{d}{dt} (U\bar{h}R^T) &= (US - U\bar{h}\bar{W}D)\bar{W}R^T \\ \tau \frac{d}{dt} \bar{h} &= \bar{W}(S - \bar{h}\bar{W}D) \end{aligned}$$

548 Here we have used the orthogonality of the singular vectors such that $V^T V = I$ and $U^T U = I$.
 549 Importantly, all matrices in the dynamics are now diagonal and represent the decoupling of the
 550 network into the modes transmitted from input to the hidden neurons and from hidden to output
 551 neurons. In practice we do not initialize the network weights to adhere to this diagonalisation and so
 552 it is not guaranteed that the matrices will be diagonal at initialisation. However, empirically it has
 553 been found that the network singular values rapidly align to this required configuration [37, 47].

554 The derivative then for the full-network input-output mapping can be obtain by using the product
 555 rule:

$$\begin{aligned} \tau \frac{d}{dt} \bar{h}\bar{W} &= (\tau \frac{d}{dt} \bar{h})\bar{W} + \bar{h}(\tau \frac{d}{dt} \bar{W}) \\ &= (\bar{W}(S - \bar{h}\bar{W}D))\bar{W} + \bar{h}(\bar{h}(S - \bar{h}\bar{W}D)) \\ &= \bar{W}^2(S - \bar{h}\bar{W}D) + \bar{h}^2(S - \bar{h}\bar{W}D) \\ &= (\bar{W}^2 + \bar{h}^2)(S - \bar{h}\bar{W}D) \end{aligned}$$

556 This means that at a minimum: $S - \overline{h}\overline{W}D = 0$ or $\frac{S}{D\overline{W}} = \overline{h}$. This defines a hyperbolic space between
 557 \overline{W} and \overline{h} . As a result we can use the change of variables: $\overline{W} = \sqrt{\lambda} \sinh \frac{\theta}{2}$ and $\overline{h} = \sqrt{\lambda} \cosh \frac{\theta}{2}$
 558 parametrized by θ .

559 We note that there is a conserved quantity between the singular values of the weight matrices:

$$\begin{aligned}\overline{W}^2 - \overline{h}^2 &= (\sqrt{\lambda} \sinh \frac{\theta}{2})^2 - (\sqrt{\lambda} \cosh \frac{\theta}{2})^2 \\ &= \lambda \sinh^2 \frac{\theta}{2} - \lambda \cosh^2 \frac{\theta}{2} \\ &= \lambda \left(\frac{\cosh(\theta) + 1}{2} \right) - \lambda \left(\frac{\cosh(\theta) - 1}{2} \right) \\ &= \frac{\lambda}{2} \cosh \theta + \frac{\lambda}{2} - \frac{\lambda}{2} \cosh \theta + \frac{\lambda}{2} \\ &= \lambda\end{aligned}$$

560 This is known as λ -Balanced weights [33] and for a given initial value for λ this quantity will be
 561 conserved for all times during training. Aiming to write the network dynamics in terms of this
 562 quantity to understand its effect on learning speed and initialisation and with the change of variables
 563 to hyperbolic coordinates we begin with:

$$\begin{aligned}(\overline{W}^2 + \overline{h}^2)^2 &= (\overline{W}^2)^2 + (\overline{h}^2)^2 \\ &= (\overline{W}^2)^2 + (\overline{h}^2)^2 + 4\overline{W}^2\overline{h}^2 - 4\overline{W}^2\overline{h}^2 \\ &= (\overline{W}^2 - \overline{h}^2)^2 + 4\overline{W}^2\overline{h}^2\end{aligned}$$

564 Substituting this into the network dynamics equation and defining the network singular value as
 565 $\omega = \overline{h}\overline{W}$ we obtain:

$$\begin{aligned}\tau \frac{d}{dt} \omega &= (\overline{W}^2 + \overline{h}^2) (S - \omega D) \\ \tau \frac{d}{dt} \omega &= \sqrt{((\overline{W}^2 - \overline{h}^2)^2 + 4\overline{W}^2\overline{h}^2)} (S - \omega D)\end{aligned}$$

566 Now applying the change of variables to hyperbolic coordinates with $\overline{W} = \sqrt{\lambda} \sinh \frac{\theta}{2}$ and $\overline{h} =$
 567 $\sqrt{\lambda} \cosh \frac{\theta}{2}$ parametrized by θ :

$$\begin{aligned}\tau \frac{d}{dt} (\sqrt{\lambda} \cosh \frac{\theta}{2}) (\sqrt{\lambda} \sinh \frac{\theta}{2}) &= \\ &= \sqrt{\left((\lambda \sinh^2 \frac{\theta}{2}) - (\lambda \cosh^2 \frac{\theta}{2}) \right)^2 + 4(\lambda \sinh^2 \frac{\theta}{2})(\lambda \cosh^2 \frac{\theta}{2})} (S - (\sqrt{\lambda} \cosh \frac{\theta}{2})(\sqrt{\lambda} \sinh \frac{\theta}{2})D) \\ \tau \frac{d}{dt} \lambda \cosh \frac{\theta}{2} \sinh \frac{\theta}{2} &= \sqrt{\left((\lambda \sinh^2 \frac{\theta}{2}) - (\lambda \cosh^2 \frac{\theta}{2}) \right)^2 + 4\lambda^2 (\cosh \frac{\theta}{2} \sinh \frac{\theta}{2})^2} (S - \lambda \cosh \frac{\theta}{2} \sinh \frac{\theta}{2} D)\end{aligned}$$

568 We can then apply the identities: $\cosh \frac{\theta}{2} \sinh \frac{\theta}{2} = \frac{1}{2} \sinh \theta$ and $\lambda \sinh^2 \frac{\theta}{2} - \lambda \cosh^2 \frac{\theta}{2} = \lambda$:

$$\begin{aligned}\tau \frac{d}{dt} \frac{\lambda}{2} \sinh(\theta) &= \sqrt{\lambda^2 + 4\lambda^2 \left(\frac{1}{2} \sinh(\theta)\right)^2} (S - \frac{\lambda}{2} \sinh(\theta)D) \\ \tau \frac{d}{dt} \frac{\lambda}{2} \sinh(\theta) &= \sqrt{\lambda^2 + \lambda^2 \sinh^2(\theta)} (S - \frac{\lambda}{2} \sinh(\theta)D) \\ \tau \frac{d}{dt} \frac{\lambda}{2} \sinh(\theta) &= |\lambda| \sqrt{1 + \sinh^2(\theta)} (S - \frac{\lambda}{2} \sinh(\theta)D)\end{aligned}$$

$$\begin{aligned}\tau \frac{d}{dt} \frac{\lambda}{2} \sinh(\theta) &= |\lambda| \sqrt{\cosh^2(\theta)} \left(S - \frac{\lambda}{2} \sinh(\theta) D\right) \\ \tau \frac{d}{dt} \frac{\lambda}{2} \sinh(\theta) &= |\lambda| \cosh(\theta) \left(S - \frac{\lambda}{2} \sinh(\theta) D\right)\end{aligned}$$

569 Now applying the derivative on the left:

$$\begin{aligned}\tau \frac{\lambda}{2} \cosh(\theta) \frac{d}{dt} \theta &= |\lambda| \cosh(\theta) \left(S - \frac{\lambda}{2} \sinh(\theta) D\right) \\ \frac{d}{dt} \theta &= \frac{1}{\tau} \operatorname{sgn}(\lambda) (2S - \lambda D \sinh(\theta))\end{aligned}$$

570 This is a separable differential equation in θ :

$$\begin{aligned}\int_{\theta_0}^{\theta_f} \frac{1}{(2S - \lambda D \sinh(\theta))} d\theta &= \int_0^t \frac{\operatorname{sgn}(\lambda)}{\tau} dt \\ \left[\frac{\log\left(\left| \frac{2S \tanh\left(\frac{\theta}{2}\right) + \sqrt{4S^2 + \lambda^2 D^2} + \lambda D \right|}{\sqrt{4S^2 + \lambda^2 D^2}}\right) - \log\left(\left| \frac{2S \tanh\left(\frac{\theta}{2}\right) - \sqrt{4S^2 + \lambda^2 D^2} + \lambda D \right|}{\sqrt{4S^2 + \lambda^2 D^2}}\right)}{\sqrt{4S^2 + \lambda^2 D^2}} \right]_{\theta_0}^{\theta_f} &= \frac{\operatorname{sgn}(\lambda)}{\tau} t \\ \frac{1}{\sqrt{4S^2 + \lambda^2 D^2}} \left[\log\left(\frac{\left| \frac{2S \tanh\left(\frac{\theta}{2}\right) + \sqrt{4S^2 + \lambda^2 D^2} + \lambda D \right|}{\left| \frac{2S \tanh\left(\frac{\theta}{2}\right) - \sqrt{4S^2 + \lambda^2 D^2} + \lambda D \right|}\right)} \right]_{\theta_0}^{\theta_f} &= \frac{\operatorname{sgn}(\lambda)}{\tau} t \\ \frac{1}{\sqrt{4S^2 + \lambda^2 D^2}} \left[\log\left(\frac{\left| \frac{2S \tanh\left(\frac{\theta_f}{2}\right) + \sqrt{4S^2 + \lambda^2 D^2} + \lambda D \right|}{\left| \frac{2S \tanh\left(\frac{\theta_f}{2}\right) - \sqrt{4S^2 + \lambda^2 D^2} + \lambda D \right|}\right)} \right. \\ \left. - \log\left(\frac{\left| \frac{2S \tanh\left(\frac{\theta_0}{2}\right) + \sqrt{4S^2 + \lambda^2 D^2} + \lambda D \right|}{\left| \frac{2S \tanh\left(\frac{\theta_0}{2}\right) - \sqrt{4S^2 + \lambda^2 D^2} + \lambda D \right|}\right)} \right] &= \frac{\operatorname{sgn}(\lambda)}{\tau} t\end{aligned}$$

571 If we let:

$$C = \frac{\left| \frac{2S \tanh\left(\frac{\theta_0}{2}\right) + \sqrt{4S^2 + \lambda^2 D^2} + \lambda D \right|}{\left| \frac{2S \tanh\left(\frac{\theta_0}{2}\right) - \sqrt{4S^2 + \lambda^2 D^2} + \lambda D \right|}; K = \sqrt{4S^2 + \lambda^2 D^2}$$

572 then:

$$\frac{1}{K} \left[\log\left(\frac{\left| \frac{2S \tanh\left(\frac{\theta_f}{2}\right) + K + \lambda D \right|}{\left| \frac{2S \tanh\left(\frac{\theta_f}{2}\right) - K + \lambda D \right|}\right)} - \log(C) \right] = \frac{\operatorname{sgn}(\lambda)}{\tau} t$$

573 Writing θ_f in terms of t :

$$\begin{aligned}\frac{1}{K} \left[\log\left(\frac{\left| \frac{2S \tanh\left(\frac{\theta_f}{2}\right) + K + \lambda D \right|}{\left| \frac{2S \tanh\left(\frac{\theta_f}{2}\right) - K + \lambda D \right|}\right)} - \log(C) \right] &= \frac{\operatorname{sgn}(\lambda)}{\tau} t \\ \log\left(\frac{\left| \frac{2S \tanh\left(\frac{\theta_f}{2}\right) + K + \lambda D \right|}{\left| \frac{2S \tanh\left(\frac{\theta_f}{2}\right) - K + \lambda D \right|}\right)} &= \frac{\operatorname{sgn}(\lambda) K}{\tau} t + \log(C) \\ \frac{\left| \frac{2S \tanh\left(\frac{\theta_f}{2}\right) + K + \lambda D \right|}{\left| \frac{2S \tanh\left(\frac{\theta_f}{2}\right) - K + \lambda D \right|}\right)} &= C \exp\left(\frac{\operatorname{sgn}(\lambda) K}{\tau} t\right) \\ 2S \tanh\left(\frac{\theta_f}{2}\right) + K + \lambda D &= C \exp\left(\frac{\operatorname{sgn}(\lambda) K}{\tau} t\right) (K - 2S \tanh\left(\frac{\theta_f}{2}\right) - \lambda D) \\ 2S \tanh\left(\frac{\theta_f}{2}\right) + C \exp\left(\frac{\operatorname{sgn}(\lambda) K}{\tau} t\right) 2S \tanh\left(\frac{\theta_f}{2}\right) &= C \exp\left(\frac{\operatorname{sgn}(\lambda) K}{\tau} t\right) (K - \lambda D) - K - \lambda D\end{aligned}$$

$$\begin{aligned}
2S \tanh\left(\frac{\theta_f}{2}\right) \left(1 + C \exp\left(\frac{\text{sgn}(\lambda)K}{\tau}t\right)\right) &= -K \left(1 - C \exp\left(\frac{\text{sgn}(\lambda)K}{\tau}t\right)\right) - \lambda D \left(1 + C \exp\left(\frac{\text{sgn}(\lambda)K}{\tau}t\right)\right) \\
\tanh\left(\frac{\theta_f}{2}\right) &= \frac{-K \left(1 - C \exp\left(\frac{\text{sgn}(\lambda)K}{\tau}t\right)\right) - \lambda D \left(1 + C \exp\left(\frac{\text{sgn}(\lambda)K}{\tau}t\right)\right)}{2S \left(1 + C \exp\left(\frac{\text{sgn}(\lambda)K}{\tau}t\right)\right)} \\
\theta_f &= 2 \tanh^{-1} \left(\frac{K \left(C \exp\left(\frac{\text{sgn}(\lambda)K}{\tau}t\right) - 1\right) - \lambda D \left(C \exp\left(\frac{\text{sgn}(\lambda)K}{\tau}t\right) + 1\right)}{2S \left(C \exp\left(\frac{\text{sgn}(\lambda)K}{\tau}t\right) + 1\right)} \right)
\end{aligned}$$

574 To obtain the dynamics for the singular value of a mode of the network we use:

$$\begin{aligned}
\omega &= \lambda \sinh \frac{\theta}{2} \cosh \frac{\theta}{2} \\
&= \frac{\lambda}{2} \sinh \theta \\
&= \frac{\lambda}{2} \sinh \left(2 \tanh^{-1} \left(\frac{K \left(C \exp\left(\frac{\text{sgn}(\lambda)K}{\tau}t\right) - 1\right) - \lambda D \left(C \exp\left(\frac{\text{sgn}(\lambda)K}{\tau}t\right) + 1\right)}{2S \left(C \exp\left(\frac{\text{sgn}(\lambda)K}{\tau}t\right) + 1\right)} \right) \right)
\end{aligned}$$

575 With the linear network dynamics we can now derive a network's hitting time (t^*). Let v^* be a
576 sufficiently small value:

$$\begin{aligned}
\frac{S}{D} - \omega &= v^* \\
\frac{S}{D} - \frac{\lambda}{2} \sinh \left(2 \tanh^{-1} \left(\frac{K \left(C \exp\left(\frac{\text{sgn}(\lambda)K}{\tau}t\right) - 1\right) - \lambda D \left(C \exp\left(\frac{\text{sgn}(\lambda)K}{\tau}t\right) + 1\right)}{2S \left(C \exp\left(\frac{\text{sgn}(\lambda)K}{\tau}t\right) + 1\right)} \right) \right) &= v^* \\
\frac{1}{2} \sinh^{-1} \left(\frac{2S - 2Dv^*}{\lambda D} \right) &= \tanh^{-1} \left(\frac{K \left(C \exp\left(\frac{\text{sgn}(\lambda)K}{\tau}t\right) - 1\right) - \lambda D \left(C \exp\left(\frac{\text{sgn}(\lambda)K}{\tau}t\right) + 1\right)}{2S \left(C \exp\left(\frac{\text{sgn}(\lambda)K}{\tau}t\right) + 1\right)} \right) \\
\tanh \left(\frac{1}{2} \sinh^{-1} \left(\frac{2S - 2Dv^*}{\lambda D} \right) \right) &= \frac{K \left(C \exp\left(\frac{\text{sgn}(\lambda)K}{\tau}t\right) - 1\right) - \lambda D \left(C \exp\left(\frac{\text{sgn}(\lambda)K}{\tau}t\right) + 1\right)}{2S \left(C \exp\left(\frac{\text{sgn}(\lambda)K}{\tau}t\right) + 1\right)}
\end{aligned}$$

577 Let $T^* = \tanh \left(\frac{1}{2} \sinh^{-1} \left(\frac{2S - 2Dv^*}{\lambda D} \right) \right)$ then

$$\begin{aligned}
T^* &= \frac{K \left(C \exp\left(\frac{\text{sgn}(\lambda)K}{\tau}t\right) - 1\right) - \lambda D \left(C \exp\left(\frac{\text{sgn}(\lambda)K}{\tau}t\right) + 1\right)}{2S \left(C \exp\left(\frac{\text{sgn}(\lambda)K}{\tau}t\right) + 1\right)} \\
2ST^* \left(C \exp\left(\frac{\text{sgn}(\lambda)K}{\tau}t\right) + 1\right) &= K \left(C \exp\left(\frac{\text{sgn}(\lambda)K}{\tau}t\right) - 1\right) - \lambda D \left(C \exp\left(\frac{\text{sgn}(\lambda)K}{\tau}t\right) + 1\right) \\
2ST^* C \exp\left(\frac{\text{sgn}(\lambda)K}{\tau}t\right) + 2ST^* &= KC \exp\left(\frac{\text{sgn}(\lambda)K}{\tau}t\right) - K - \lambda DC \exp\left(\frac{\text{sgn}(\lambda)K}{\tau}t\right) - \lambda D \\
2ST^* C \exp\left(\frac{\text{sgn}(\lambda)K}{\tau}t\right) - KC \exp\left(\frac{\text{sgn}(\lambda)K}{\tau}t\right) + \lambda DC \exp\left(\frac{\text{sgn}(\lambda)K}{\tau}t\right) &= -2ST^* - K - \lambda D \\
\exp\left(\frac{\text{sgn}(\lambda)K}{\tau}t\right) (2ST^* C - KC + \lambda DC) &= -2ST^* - K - \lambda D \\
\exp\left(\frac{\text{sgn}(\lambda)K}{\tau}t\right) &= \frac{-2ST^* - K - \lambda D}{2ST^* C - KC + \lambda DC}
\end{aligned}$$

$$\begin{aligned}\frac{\operatorname{sgn}(\lambda)K}{\tau}t &= \log\left(\frac{-2ST^* - K - \lambda D}{2ST^*C - KC + \lambda DC}\right) \\ t^* &= \frac{\tau}{\operatorname{sgn}(\lambda)K} \log\left(\frac{-2ST^* - K - \lambda D}{2ST^*C - KC + \lambda DC}\right) \\ t^* &= \frac{\tau}{\operatorname{sgn}(\lambda)K} \log\left(\frac{K + 2ST^* + \lambda D}{KC - 2ST^*C - \lambda DC}\right)\end{aligned}$$

578 Similarly we derive the escaping time for a mode with sufficiently small \hat{v} as:

$$\begin{aligned}\omega &= \hat{v} \\ \frac{\lambda}{2} \sinh\left(2 \tanh^{-1}\left(\frac{K\left(C \exp\left(\frac{\operatorname{sgn}(\lambda)K}{\tau}t\right) - 1\right) - \lambda D\left(C \exp\left(\frac{\operatorname{sgn}(\lambda)K}{\tau}t\right) + 1\right)}{2S\left(C \exp\left(\frac{\operatorname{sgn}(\lambda)K}{\tau}t\right) + 1\right)}\right)\right) &= \hat{v} \\ \frac{K\left(C \exp\left(\frac{\operatorname{sgn}(\lambda)K}{\tau}t\right) - 1\right) - \lambda D\left(C \exp\left(\frac{\operatorname{sgn}(\lambda)K}{\tau}t\right) + 1\right)}{2S\left(C \exp\left(\frac{\operatorname{sgn}(\lambda)K}{\tau}t\right) + 1\right)} &= \tanh\left(\frac{1}{2} \sinh^{-1}\left(\frac{2\hat{v}}{\lambda}\right)\right)\end{aligned}$$

579 Let $\hat{T} = \tanh\left(\frac{1}{2} \sinh^{-1}\left(\frac{2\hat{v}}{\lambda}\right)\right)$ then

$$\begin{aligned}\hat{T} &= \frac{K\left(C \exp\left(\frac{\operatorname{sgn}(\lambda)K}{\tau}t\right) - 1\right) - \lambda D\left(C \exp\left(\frac{\operatorname{sgn}(\lambda)K}{\tau}t\right) + 1\right)}{2S\left(C \exp\left(\frac{\operatorname{sgn}(\lambda)K}{\tau}t\right) + 1\right)} \\ 2S\hat{T}\left(C \exp\left(\frac{\operatorname{sgn}(\lambda)K}{\tau}t\right) + 1\right) &= K\left(C \exp\left(\frac{\operatorname{sgn}(\lambda)K}{\tau}t\right) - 1\right) - \lambda D\left(C \exp\left(\frac{\operatorname{sgn}(\lambda)K}{\tau}t\right) + 1\right) \\ 2S\hat{T}C \exp\left(\frac{\operatorname{sgn}(\lambda)K}{\tau}t\right) + 2S\hat{T} &= KC \exp\left(\frac{\operatorname{sgn}(\lambda)K}{\tau}t\right) - K - \lambda DC \exp\left(\frac{\operatorname{sgn}(\lambda)K}{\tau}t\right) - \lambda D \\ 2S\hat{T}C \exp\left(\frac{\operatorname{sgn}(\lambda)K}{\tau}t\right) - KC \exp\left(\frac{\operatorname{sgn}(\lambda)K}{\tau}t\right) &+ \lambda DC \exp\left(\frac{\operatorname{sgn}(\lambda)K}{\tau}t\right) = -2S\hat{T} - K - \lambda D \\ \exp\left(\frac{\operatorname{sgn}(\lambda)K}{\tau}t\right) (2S\hat{T}C - KC + \lambda DC) &= -2S\hat{T} - K - \lambda D \\ \exp\left(\frac{\operatorname{sgn}(\lambda)K}{\tau}t\right) &= \frac{-2S\hat{T} - K - \lambda D}{2S\hat{T}C - KC + \lambda DC} \\ \frac{\operatorname{sgn}(\lambda)K}{\tau}t &= \log\left(\frac{-2S\hat{T} - K - \lambda D}{2S\hat{T}C - KC + \lambda DC}\right) \\ \hat{t} &= \frac{\tau}{\operatorname{sgn}(\lambda)K} \log\left(\frac{-2S\hat{T} - K - \lambda D}{2S\hat{T}C - KC + \lambda DC}\right) \\ \hat{t} &= \frac{\tau}{\operatorname{sgn}(\lambda)K} \log\left(\frac{K + 2S\hat{T} + \lambda D}{KC - 2S\hat{T}C - \lambda DC}\right)\end{aligned}$$

580 Thus, the escaping time can be summarised as:

$$\hat{t} = \frac{\tau}{\operatorname{sgn}(\lambda)K} \log\left(\frac{K + 2S\hat{T} + \lambda D}{KC - 2S\hat{T}C - \lambda DC}\right) \quad (13)$$

581 with the escaping time constant:

$$\hat{T} = \tanh\left(\frac{1}{2} \sinh^{-1}\left(\frac{2\hat{v}}{\lambda}\right)\right) \quad (14)$$

582 Similarly the hitting time is summarised as:

$$t^* = \frac{\tau}{\text{sgn}(\lambda)K} \log \left(\frac{K + 2ST^* + \lambda D}{KC - 2ST^*C - \lambda DC} \right) \quad (15)$$

583 with the hitting time constant:

$$T^* = \tanh \left(\frac{1}{2} \sinh^{-1} \left(\frac{2S - 2Dv^*}{\lambda D} \right) \right) \quad (16)$$

584 B Mean-field theory of the dynamics

585 As outlined in Sec.4, the key observation for the mean-field analysis is that the main properties of
 586 the learning dynamics can be expressed as functions of the order parameters—Eqs. 8. By combining
 587 these definitions with the update rules—Eqs. (6, 7)—we can derive closed-form expressions for the
 588 evolution of the order parameters, enabling us to track the key observables throughout the training
 589 process. In the high-dimensional limit ($d \rightarrow \infty$), these discrete update equations converge to ordinary
 590 differential equations (ODEs), which can be integrated either numerically or analytically in certain
 591 cases [77]. As is often the case in the statistical physics of disordered systems, this approach was first
 592 derived non-rigorously by [38] and [40], with later works laying down a mathematical foundation
 593 showing concentration of the ODEs [73, 78].

594 Following these prescriptions, we obtain the update equations as in [42]. Let us define the pre-
 595 activations of the student and task- t teacher given an input \mathbf{x} from task t as

$$\lambda_i = \frac{1}{\sqrt{d}} \mathbf{W}_i \cdot \mathbf{x}, \quad \rho_i^{(t)} = \frac{1}{\sqrt{d}} \mathbf{W}_{T,i}^{(t)} \cdot \mathbf{x}, \quad (17)$$

596 and denote the difference between the teacher and student predictions by $\Delta^{(t)} = \mathbf{h}^{(t)} \cdot \phi(\boldsymbol{\lambda}) - \mathbf{h}_T^{(t)} \cdot \phi(\boldsymbol{\rho})$.
 597 The corresponding ODEs for the order parameters in the limit $d \rightarrow \infty$ are given by:

$$\frac{dQ_{ik}}{d\tau} = -\eta h_i^{(t)} \langle \phi'(\lambda_i) \Delta^{(t)} \lambda_k \rangle - \eta h_k^{(t)} \langle \phi'(\lambda_k) \Delta^{(t)} \lambda_i \rangle + \eta^2 h_i^{(t)} h_k^{(t)} \langle \phi'(\lambda_i) \phi'(\lambda_k) (\Delta^{(t)})^2 \rangle, \quad (18)$$

$$\frac{dR_{in}^{(t')}}{d\tau} = -\eta h_i^{(t)} \langle \phi'(\lambda_i) \Delta^{(t)} \rho_n^{(t')} \rangle, \quad (19)$$

$$\frac{dh_i^{(t)}}{d\tau} = -\eta \langle \Delta^{(t)} \phi(\lambda_i) \rangle, \quad (20)$$

598 where $\tau = \text{epoch}/d$ represents continuous time in the high-dimensional limit, and $t, t' \in 1, 2$
 599 denote the task indices. The angular brackets indicate an average over the pre-activations. The
 600 pre-activations themselves are centered Gaussian random variables with covariances determined by
 601 the order parameters \mathbf{Q} , $\mathbf{R}^{(t)}$, and \mathbf{T} .

602 These averages can be computed analytically for certain activation functions. For instance, in the case
 603 of a rescaled error function introduced in the main text [38, 40], the relevant averages are given by:

$$\langle \phi(\beta) \phi(\gamma) \rangle = \frac{1}{\pi} \arcsin \left(\frac{\Sigma_{12}}{\sqrt{(1 + \Sigma_{11})(1 + \Sigma_{22})}} \right), \quad (21)$$

$$\langle \phi'(\zeta) \beta \phi(\gamma) \rangle = \frac{2\Sigma_{23}(1 + \Sigma_{11}) - 2\Sigma_{12}\Sigma_{13}}{\sqrt{\Lambda_3}(1 + \Sigma_{11})}, \quad (22)$$

$$\langle \phi'(\zeta) \phi'(\iota) \phi(\beta) \phi(\gamma) \rangle = \frac{4}{\pi^2 \sqrt{\Lambda_4}} \arcsin \left(\frac{\Lambda_0}{\sqrt{\Lambda_1 \Lambda_2}} \right), \quad (23)$$

604 where the Greek letters represent arbitrary pre-activations with covariance matrix $\boldsymbol{\Sigma}$, and the auxiliary
 605 quantities Λ_i are given by:

$$\Lambda_0 = \Lambda_4 \Sigma_{34} - \Sigma_{23} \Sigma_{24} (1 + \Sigma_{11}) - \Sigma_{13} \Sigma_{14} (1 + \Sigma_{22}) + \Sigma_{12} \Sigma_{13} \Sigma_{24} + \Sigma_{12} \Sigma_{14} \Sigma_{23}, \quad (24)$$

$$\Lambda_1 = \Lambda_4 (1 + \Sigma_{33}) - \Sigma_{23}^2 (1 + \Sigma_{11}) - \Sigma_{13}^2 (1 + \Sigma_{22}) + 2\Sigma_{12} \Sigma_{13} \Sigma_{23}, \quad (25)$$

$$\Lambda_2 = \Lambda_4 (1 + \Sigma_{44}) - \Sigma_{24}^2 (1 + \Sigma_{11}) - \Sigma_{14}^2 (1 + \Sigma_{22}) + 2\Sigma_{12} \Sigma_{14} \Sigma_{24}, \quad (26)$$

$$\Lambda_3 = (1 + \Sigma_{11})(1 + \Sigma_{33}) - \Sigma_{13}^2. \quad (27)$$

606 These expressions provide a comprehensive analytical framework for tracking the dynamics of the
 607 student network and the evolution of specialisation across training.

608 C Method for Linear Network Phase Transition

609 D Disentanglement

610 We conduct our experiments using open-source frameworks [24, 79]. Specifically, we implement
611 a beta-VAE with the "DeepGaussianLinear" architecture for the decoder and "DeepLinear" for the
612 encoder. We modify the Xavier initialisation where the weights of the linear layers will have values
613 sampled from $U(-a, a)$ with

$$a = \text{gain} \times \sqrt{\frac{6}{\text{fan_in} + \text{fan_out}}}$$

614 We vary the gain between 0.3 and 3 and run each experiment over 4 seeds. All network parameters
615 are set to their default values as provided by the respective open-source frameworks. We run the
616 experiments for 20 Epochs and 157499 iterations.

617 These experiment illustrate the impact of initialisation on network specialisation. Although the scope
618 of these experiments is limited, they provide preliminary validation of our theoretical framework in
619 more realistic contexts. We advocate for further investigation into alternative initialisation schemes
620 with varying levels of balance. Moreover, we highlight the need for future research to extend
621 these experiments by considering a wider variety of datasets (Car3D [80], dSprites [76],), network
622 architectures (Conv,Linear), initialisation strategies (Gaussian Xavier Initalisation) and different
623 metric (SAP [68, 63],) to fully explore the implications of our findings. In practice, linear networks
624 in PyTorch are initialized using a uniform distribution, specifically:

$$\mathbf{W} \sim \mathcal{U}(-\sqrt{k}, \sqrt{k}), \quad \text{where } k = \frac{1}{\text{in_features}}$$

625 This initialization is equivalent to applying a small gain in our experimental setting, aligning with the
626 weight scaling typically seen in neural network training setups. **DCI Disentanglement** [72] define
627 three key properties of learned representations: Disentanglement, Completeness, and Informativeness.
628 To assess these, they calculate the importance of each dimension of the representation in predicting
629 a factor of variation. This can be done using models like Lasso or Random Forest classifiers.
630 Disentanglement is computed by subtracting the entropy of the probability that a representation
631 dimension predicts a factor, weighted by its relative importance. Completeness is similarly measured,
632 focusing on how well a factor is captured by the dimensions. Informativeness is evaluated as the
633 prediction error of the factors. We use the implementation in [24]. In this implementation, we sample
634 10,000 training and 5,000 test points, then use gradient-boosted trees from Scikit-learn to obtain
635 feature importance weights. These weights form an importance matrix, with rows representing factors
636 and columns representing dimensions. Disentanglement is calculated by normalizing the columns of
637 this matrix, subtracting the entropy from 1 for each column, and then weighting by each dimension's
638 relative importance.

639 E Additional Entropy Phase Diagrams

640 In Fig. 5 we showed phase diagrams of the aggregate entropy as a function of initialisation parameters,
641 for both ReLU and sigmoidal networks. Below we show additional plots with the individual entropy
642 terms (H_u defined over the unit activations, and H_h defined over the head weights).

643 F Diversity of Forgetting Curves

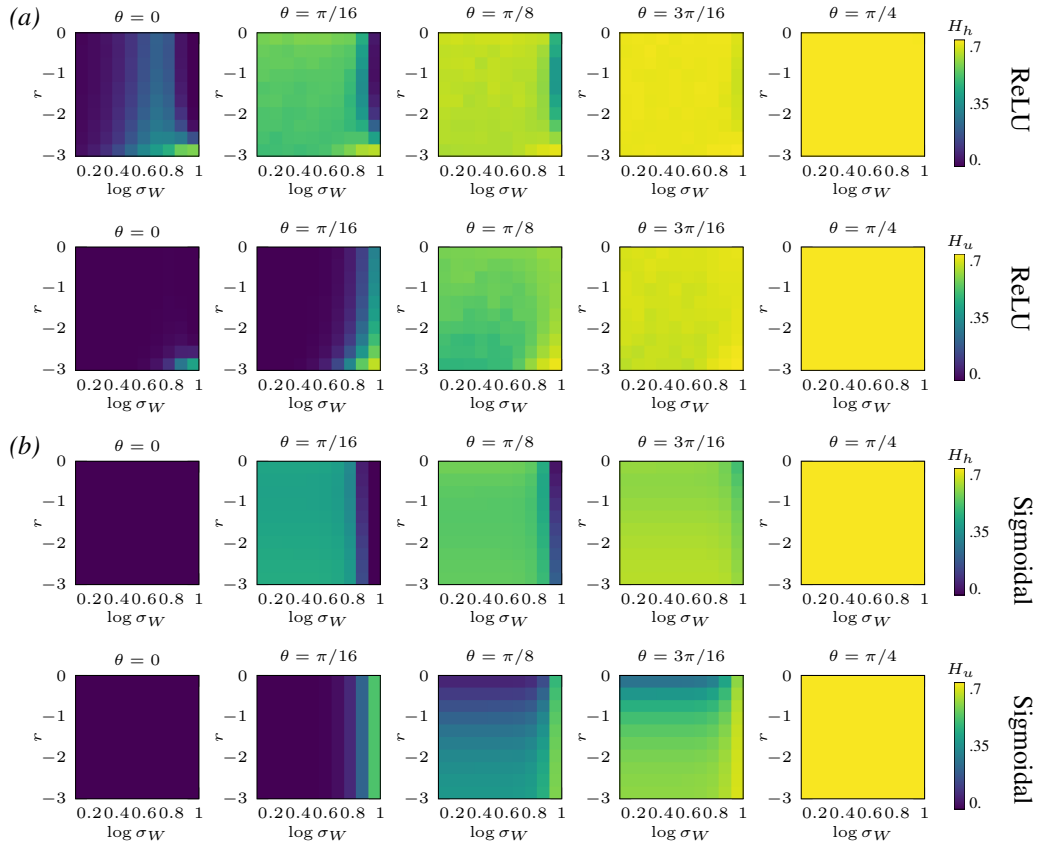


Figure 8: **Additional Phase Diagrams.** Here we show the equivalent phase diagrams from Fig. 5 for entropy measures over the unit activations and head weights.

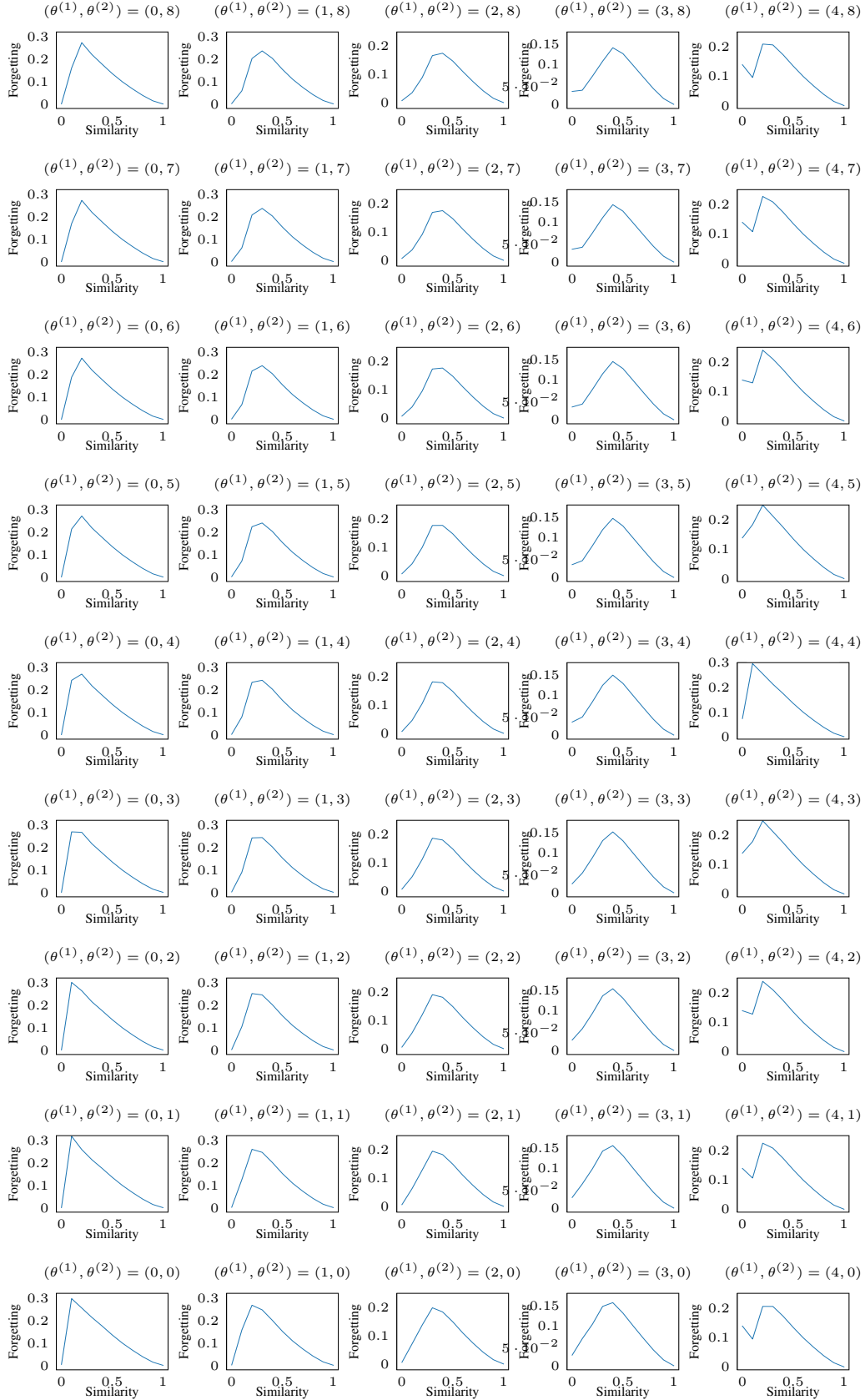


Figure 9: **Initialisation can lead to a diversity of specialisation dynamics and a diversity of relationships between forgetting and task similarity.** R, σ_W fixed, $\theta^{(1)}, \theta^{(2)}$ measured in increments of $\pi/16$. Scaled error function, $P^* = 1, P = 1$. 24