

---

# EARLY LEARNING OF THE OPTIMAL CONSTANT SOLUTION IN NEURAL NETWORKS AND HUMANS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Deep neural networks learn increasingly complex functions over the course of training. Here, we show both empirically and theoretically that learning of the target function is preceded by an early phase in which networks learn the optimal constant solution (OCS) – that is, initial model responses mirror the distribution of target labels, while entirely ignoring information provided in the input. Using a hierarchical category learning task, we derive exact solutions for learning dynamics in deep linear networks trained with bias terms. Even when initialized to zero, this simple architectural feature induces substantial changes in early dynamics. We identify hallmarks of this early OCS phase and illustrate how these signatures are observed in deep linear networks and larger, more complex (and nonlinear) convolutional neural networks solving a hierarchical learning task based on MNIST and CIFAR10. We explain these observations by proving that deep linear networks necessarily learn the OCS during early learning. To further probe the generality of our results, we train human learners over the course of three days on a structurally equivalent learning task. We then identify qualitative signatures of this early OCS phase in terms of true negative rates. Surprisingly, we find the same early reliance on the OCS in the behaviour of human learners. Finally, we show that learning of the OCS can emerge even in the absence of bias terms and is equivalently driven by generic correlations in the input data. Overall, our work suggests the OCS as a common learning principle in supervised, error-corrective learning, and suggests possible factors for its prevalence.

## 1 INTRODUCTION

Neural networks trained with stochastic gradient descent (SGD) exhibit various *simplicity biases*, where models tend to learn simple functions before more complex ones (Kalimeris et al., 2019; Rahaman et al., 2019). Simplicity biases hold significant theoretical interest as they provide an explanation for how deep networks generalize or fail to generalize in practice (Bhattacharjee et al., 2023; Valle-Pérez et al., 2019; Zhang et al., 2021).

The characterisation of simplicity biases is still incomplete. Some explanations appeal to distributional properties of input data, pointing out that SGD progressively learns increasingly higher-order moments (Refinetti et al., 2023; Belrose et al., 2024). Other approaches focus directly on the evolution of the network function, proposing that networks initially learn a classifier highly correlated with a linear model. Importantly, networks continue to perform well on examples correctly classified by this simple function, even when overfitting in later training (Kalimeris et al., 2019). This implies that dynamical simplicity biases help models generalize, by locking in initial knowledge that is not erased or forgotten during later training (Braun et al., 2022; Kalimeris et al., 2019).

Deep linear networks have proven to be a valuable tool for studying simplicity biases. A key finding is that directions in the network function are learned in order of importance (Saxe et al., 2014; 2019). This phenomenon, known as *progressive differentiation*, connects modern deep learning theory to both to human child development and to the earliest connectionist models of semantic cognition (Rogers and McClelland, 2004; Rumelhart et al., 1986).

Our contribution proposes a connection between these works by characterizing networks in the *earliest* stages of learning in terms of input, output, and architecture. In the hierarchical setting by Saxe et al. (2019), we demonstrate both theoretically and empirically that neural networks initially

054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082  
083  
084  
085  
086  
087  
088  
089  
090  
091  
092  
093  
094  
095  
096  
097  
098  
099  
100  
101  
102  
103  
104  
105  
106  
107

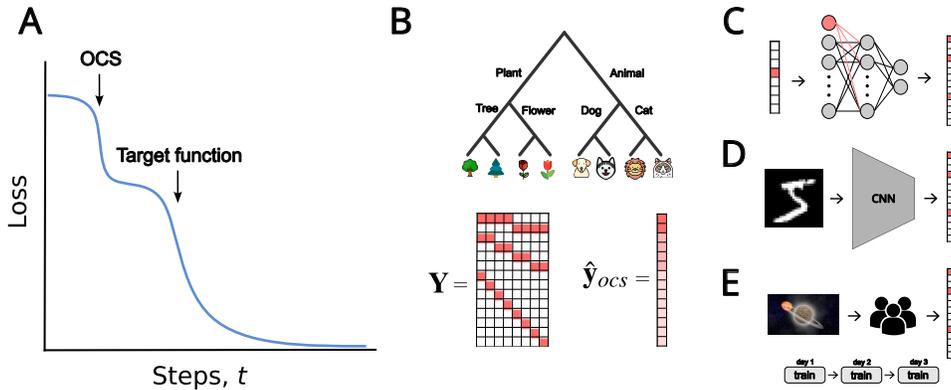


Figure 1: Early learning of the optimal constant solution (OCS). **A** Graphical illustration of our hypothesis, where learning of the target function is preceded by the early acquisition of the OCS. **B Top**: a graphical illustration of the hierarchical structure embedded in the outputs. **Bottom**: The full output data matrix  $\mathbf{Y}$  used across different types of learners and the corresponding OCS solution  $\hat{\mathbf{y}}_{ocs}$ . **C** Illustration of experiments in linear networks with bias terms. **D** Illustration of our experiments with non-linear models and **E** Illustration of the task as adapted for humans.

learn via the output statistics of the data. This function has been termed the optimal constant solution (OCS) by Kang et al. (2024), who demonstrated that networks revert to the OCS when probed on out-of-distribution inputs. Here, we demonstrate and prove how linear networks, when equipped with these bias terms, necessarily learn the OCS early in training. Fig. 1A graphically illustrates this observation. We furthermore highlight the practical relevance of these results by examining early learning dynamics in complex, non-linear architectures.

Biological learners also display behaviours that imply the input-independent learning of output statistics. In probability matching, responses mirror the probabilities of rewarded actions (Herrnstein, 1961; Estes, 1964; Estes and Straughan, 1954). Learners often display non-stationary biases that are driven by the distribution of recent responses (Jones et al., 2015; Gold et al., 2008; Verplanck et al., 1952). In paired-associates learning accuracy can depend not only on a learned input-output mapping but also on knowledge of the task structure (Hawker, 1964; Bower, 1962). Humans also display simplicity biases and preferentially use simple over complex functions (Feldman, 2000; Goodman et al., 2008; Chater, 1996; Lombrozo, 2007; Feldman, 2003). However, relatively little attention has been devoted to the dynamics of these biases. We conduct experiments to determine whether humans replicate early reliance on the OCS.

### 1.1 CONTRIBUTIONS

- We devise exact solutions for learning dynamics to analyse linear networks with bias in the input layer. Even when initialized at zero, this component substantially alters *early* learning dynamics.
- We empirically characterise early learning in these linear networks as being dominated by average output statistics. We explain this result with a theoretical analysis which reveals that average output statistics are always learned first when the network contains bias terms.
- We further highlight the practical relevance of these theoretical results in a hierarchical learning task for humans as well as complex, non-linear architectures by empirically demonstrating that all learners develop stereotypical response biases during early stages of training.
- On the basis of the developed theory we show that, in linear networks, early OCS learning can be induced by input correlations even in absence of bias terms. For natural datasets we empirically demonstrate that learning of the OCS can indeed be purely driven by generic correlations in the input data.

---

## 1.2 RELATED WORK

**Deep linear networks.** In deep linear networks analytical solutions have been obtained for certain initial conditions and datasets (Saxe et al., 2014; 2019; Braun et al., 2022; Fukumizu, 1998). Progress has also been made in understanding linear network loss landscapes (Baldi and Hornik, 1989) and generalisation ability (Lampinen and Ganguli, 2019). Despite their linearity these models display complex non-linear learning dynamics which reflect behaviours seen in non-linear models (Saxe et al., 2019). Moreover, learning dynamics in such simple models have been argued to qualitatively resemble phenomena observed in the cognitive development of humans (Saxe et al., 2019; Rogers and McClelland, 2004).

**Biological response biases.** Humans and animals routinely display response biases during perceptual learning and decision making tasks (Gold et al., 2008; Jones et al., 2015; Liebana Garcia et al., 2023; Amitay et al., 2014; Urai et al., 2019). In these tasks decisions are frequently made in sequences where responses and feedback steer decisions beyond the provided perceptual evidence (Jones et al., 2015; Fan et al., 2024; Gold et al., 2008; Verplanck et al., 1952; Sugrue et al., 2004). Non-stationary response biases can be driven by feedback on previous trials (Dutilh et al., 2012; Rabbitt and Rodgers, 1977) or might reflect global beliefs about the statistics of a task (Fan et al., 2024; Jones et al., 2015). Importantly, response biases are particularly pronounced in early learning (Jones et al., 2015; Gold et al., 2008; Liebana Garcia et al., 2023) and their influence appears to be strongest when uncertainty about the correct response is highest (Gold et al., 2008; Fan et al., 2024).

**Simplicity biases in machine learning.** Simplicity biases in neural networks have been studied extensively both theoretically (Bordelon et al., 2020; Mei et al., 2022) and empirically (Bhattamishra et al., 2023; Mingard et al., 2023). Work on the *distributional* simplicity bias emphasises the importance of input data and proposes that models learn via progressive exploitation of dataset moments (Refinetti et al., 2023; Belrose et al., 2024). On the other hand, neural networks have been found to express simpler functions during early training (Kalimeris et al., 2019; Refinetti et al., 2023; Belrose et al., 2024; Rahaman et al., 2019). Our work draws a connection between these findings and highlights how input statistics bias early learning towards output statistics.

## 1.3 PAPER ORGANISATION

We initially review the linear network formalism in Section 2 on which we base our theoretical analysis. In Section 3 we derive learning dynamics for linear networks with bias terms trained on a classic hierarchical task and we document substantial changes in early dynamics. Section 4 characterizes this period of early learning empirically, and provides a theoretical explanation. We then in Section 4 validate the relevance of our findings for learning in complex models. Section 5 demonstrates the prevalence of early OCS learning in humans. Finally, Section 6 further probes generality by considering natural datasets and models that do not strictly fulfil the previous theoretical assumptions.

## 2 LINEAR NETWORK PRELIMINARIES

Here, we briefly review the analytical approach to learning dynamics in linear networks developed by Saxe et al. (2014; 2019). Consider a learning task in which a network is presented with input vectors  $\mathbf{x}_i \in \mathbb{R}^{N_{in}}$  that are associated to output vectors  $\mathbf{y}_i \in \mathbb{R}^{N_{out}}$ . The total dataset consists of  $\{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^N$  with  $N$  samples. For our setting we consider two layer linear networks where the forward pass computes  $\hat{\mathbf{y}}_i = \mathbf{W}^2 \mathbf{W}^1 \mathbf{x}_i$  and shallow networks with forward pass  $\hat{\mathbf{y}}_i = \mathbf{W}^s \mathbf{x}_i$ . Here weight matrices are of dimension  $\mathbf{W}^1 \in \mathbb{R}^{N_{hid} \times N_{in}}$ ,  $\mathbf{W}^2 \in \mathbb{R}^{N_{out} \times N_{hid}}$ , and  $\mathbf{W}^s \in \mathbb{R}^{N_{out} \times N_{in}}$ . We train our networks to minimise a squared error loss of the form  $\mathcal{L}(\hat{\mathbf{y}}) = \frac{1}{2} \sum_{i=1}^N \|\mathbf{y}_i - \hat{\mathbf{y}}_i\|^2$ .

We optimise networks using full batch-gradient descent in the gradient flow regime. When learning from small initial conditions, dynamics in these simple networks are solely dependent on the dataset input-output and input-input correlation matrices (Saxe et al., 2014). Using singular value decomposition (SVD), these matrices can be expressed as

$$\Sigma^{yx} = \frac{1}{N} \mathbf{YX}^T = \mathbf{USV}^T, \quad \Sigma^x = \frac{1}{N} \mathbf{XX}^T = \mathbf{VDV}^T. \quad (1)$$

162  
163  
164  
165  
166  
167  
168  
169  
170  
171  
172  
173  
174  
175  
176  
177  
178  
179  
180  
181  
182  
183  
184  
185  
186  
187  
188  
189  
190  
191  
192  
193  
194  
195  
196  
197  
198  
199  
200  
201  
202  
203  
204  
205  
206  
207  
208  
209  
210  
211  
212  
213  
214  
215

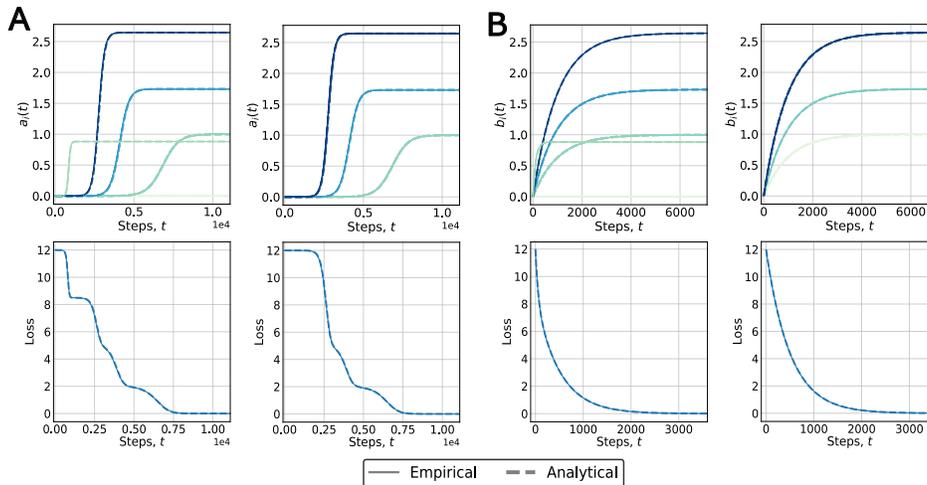


Figure 2: Exact learning dynamics. **A** Deep linear networks with bias (left) and without bias term (right). **B** Shallow linear networks with bias (left) and without bias term (right). *Top row*: Exact and simulated effective singular values  $\mathbf{A}(t)$  and  $\mathbf{B}(t)$  for deep and shallow linear networks respectively. Different  $a_\alpha(t)$  and  $b_\alpha(t)$  are color-coded according to their asymptote value with larger values as darker. *Bottom row*: Exact and simulated loss.

Here  $\mathbf{X} \in \mathbb{R}^{N_{in} \times N}$  and  $\mathbf{Y} \in \mathbb{R}^{N_{out} \times N}$  contain the full set of input vectors and output vectors. Crucially, if the right singular vectors  $\mathbf{V}^T$  of  $\Sigma^{yx}$  diagonalise  $\Sigma^x$  (see Proposition 1) then the full evolution of network weights for deep and shallow networks through time can be described as

$$\mathbf{W}^2(t)\mathbf{W}^1(t) = \mathbf{U}\mathbf{A}(t)\mathbf{V}^T. \quad (2)$$

Here  $\mathbf{A}(t)$  is a diagonal matrix. The evolution of these diagonal values  $\mathbf{A}(t)_{\alpha\alpha} = a_\alpha(t)$  at each time-step  $t$  then follows a sigmoidal trajectory as expressed in Eq. (3). For shallow networks we can similarly describe the evolution of the weight matrix  $\mathbf{W}^s(t)$  as  $\mathbf{U}\mathbf{B}(t)\mathbf{V}^T$ . Here the diagonal values  $\mathbf{B}(t)_{\alpha\alpha} = b_\alpha(t)$  evolve as seen in Eq. (4)

$$a_\alpha(t) = \frac{s_\alpha/d_\alpha}{1 - (1 - \frac{s_\alpha}{d_\alpha a_0})e^{-\frac{2s_\alpha}{\tau}t}} \quad (3) \quad b_\alpha(t) = \frac{s_\alpha}{d_\alpha}(1 - e^{-\frac{d_\alpha}{\tau}t}) + b_0 e^{-\frac{d_\alpha}{\tau}t} \quad (4)$$

In Eq. (3)  $s_\alpha = \mathbf{S}_{\alpha\alpha}$  and  $d_\alpha = \mathbf{D}_{\alpha\alpha}$  denote the relevant singular values of  $\Sigma^{yx}$  and the eigenvalues of  $\Sigma^x$  respectively,  $a_0$  are the singular values at initialisation, and  $\tau = \frac{1}{N\epsilon}$  is the time constant where  $\epsilon$  is the learning rate. In Eq. (4)  $b_0$  is the initial condition given by the initialisation. Importantly, these relations reveal that singular values control learning speed. These solutions hinge on the diagonalisation of  $\Sigma^x$  through  $\mathbf{V}$ . Prior work has focused on the case of white inputs, i.e.  $\Sigma^x = \mathbf{I}_N$  where  $\mathbf{I}_N$  denotes the  $N \times N$  identity matrix. The solution holds trivially as any  $\mathbf{V}$  will orthogonalise  $\Sigma^x$  (Saxe et al., 2019). We discuss a relevant relaxation of this condition in Proposition 1. While solutions can be derived for some non-white inputs, little attention has been devoted to learning dynamics in these scenarios. We will show how these solutions apply when networks contain bias terms in the input layer.

### 3 EXACT LEARNING DYNAMICS WITH BIAS TERMS

In this section, we derive exact learning dynamics in linear networks with bias terms and analyse the resulting changes in the dynamics. This extension to the theory by Saxe et al. (2014) forms the basis for our later discussion. For simplicity, we focus on input bias terms and uncorrelated data, but explore bias terms in other layers and correlated inputs in Appendix A.5.2 and Section 6, respectively.

**Closed-form learning dynamics.** We consider uncorrelated inputs  $\mathbf{X} = \mathbf{I}_N$  where  $\mathbf{I}_N$  denotes the  $N \times N$  identity matrix. Our linear network with a bias term in the input layer will compute  $\mathbf{W}^2(\tilde{\mathbf{W}}^1 \mathbf{x}_i + \tilde{\mathbf{b}})$  where  $\tilde{\mathbf{b}}$  are learnable bias terms.

A priori, it is unclear whether the diagonalisation of  $\Sigma^x$  through  $\mathbf{V}$  in Eq. (1) is possible in presence of bias terms. Here, we state the condition under which learning dynamics can be described in closed-form.

**Proposition 1** (Feasibility of closed-form learning dynamics). *For any input data  $\mathbf{X} \in \mathbb{R}^{N_{in} \times N}$  and output data  $\mathbf{Y} \in \mathbb{R}^{N_{out} \times N}$  it is possible to diagonalize  $\Sigma^x$  by the right singular vectors  $\mathbf{V}$  of  $\Sigma^{y^x}$  if  $\mathbf{Y}^T \mathbf{Y}$  and  $\mathbf{X}^T \mathbf{X}$  commute. The converse holds true only if  $\mathbf{X}$  has a left inverse.*

A proof is given in Appendix A.5.4. We put this statement to use to assess the effect of a bias term on learning, building on the formalism from Section 2. To this end, we re-express the network weights as  $\mathbf{W}^1 = [\tilde{\mathbf{b}} \quad \tilde{\mathbf{W}}^1]$  with inputs defined as  $\mathbf{x}_i = [1 \quad \mathbf{I}_i^T]^T$  where  $\mathbf{I}_i$  denotes the  $i$ th column of the  $N \times N$  identity matrix (see Appendix A.5.1). To introduce a controlled setting in which to analyze the effect of bias terms, we now first consider a canonical hierarchical learning task while later sections of the paper will generalize our findings beyond this setting.

**The hierarchical task.** The hierarchical task requires learning a mapping from one-hot, input vectors to output vectors that are depicted in Fig. 1B. Hereby each output vector is “three-hot”, i.e. the vector has three entries/labels. The hierarchical structure arises from the similarity between output vectors where some labels  $y^m(\mathbf{x}_i)$  are more general and correspond to more than one input  $\mathbf{x}_i$ , while labels corresponding to the bottom of the hierarchy are specific to a single input vector  $\mathbf{x}_i$ . The task is motivated in the literature on semantic cognition and leverages the fact that semantic information is usually hierarchically structured (Rogers and McClelland, 2004). In Fig. 2 we depict exact learning trajectories for the hierarchically structured outputs from Fig. 1B.

Importantly, the introduction of a bias term  $\mathbf{X} \rightarrow [\mathbf{1}_N \quad \mathbf{X}]^T$  does not affect the commutativity of  $\mathbf{X}^T \mathbf{X}$  and  $\mathbf{Y}^T \mathbf{Y}$  for the hierarchical dataset, as the constant mode  $\mathbf{1}_N$  (i.e., a vector of 1s) is already an eigenvector to both these similarity matrices (see Appendix A.5.6). In consequence, the analytical solutions in Section 2 remain applicable. We generalize these considerations in Section 6 and Appendix A.5.6.

Fig. 2 shows that linear networks with bias terms have a distinctly different early learning phase when compared to vanilla linear networks. While both models converge to a zero loss solution, we observe that the final network function with bias terms contains an additional non-zero singular value with their associated singular vectors. We devote the next section to analyzing this change in the early dynamics.

## 4 BIAS TERMS DRIVE EARLY LEARNING TOWARDS THE OPTIMAL CONSTANT SOLUTION

In this section, we qualitatively characterize what causes observed changes in early learning dynamics. We find that early learning dynamics are driven by average output statistics and provide a theoretical explanation. We then demonstrate the generality of this result by highlighting how early learning of average output statistics can be similarly observed in complex, non-linear architectures.

A naive strategy to learning is to minimise error over a set of samples while disregarding information conveyed by the input. Previous work has recently termed this network function the optimal constant solution (OCS) (Kang et al., 2024). The OCS can be formalised as  $\hat{\mathbf{y}}_{ocs} = \operatorname{argmin}_{\hat{\mathbf{y}} \in \mathbb{R}^{N_{out}}} \frac{1}{N} \sum_{i=1}^N \mathcal{L}(\hat{\mathbf{y}}, \mathbf{y}_i)$  and represents the optimal function  $\hat{\mathbf{y}}$  that is independent of input  $\mathbf{x}_i$ . For mean-squared error, it is straightforward to show that the minimiser is the average output  $\hat{\mathbf{y}}_{ocs} = \frac{1}{N} \sum_i \mathbf{y}_i =: \bar{\mathbf{y}}$ .

**Setup.** We train linear networks and Convolutional neural networks (CNN) on the hierarchical learning task illustrated in Fig. 1C and D respectively. For CNNs we design a “hierarchical MNIST” task whereby one-hot inputs are replaced with eight randomly sampled classes from MNIST (Li Deng, 2012). For the “hierarchical MNIST” task we used the started from ten digit classes provided by MNIST and then sampled 8 classes randomly. For each image in each class we then replaced

270  
271  
272  
273  
274  
275  
276  
277  
278  
279  
280  
281  
282  
283  
284  
285  
286  
287  
288  
289  
290  
291  
292  
293  
294  
295  
296  
297  
298  
299  
300  
301  
302  
303  
304  
305  
306  
307  
308  
309  
310  
311  
312  
313  
314  
315  
316  
317  
318  
319  
320  
321  
322  
323

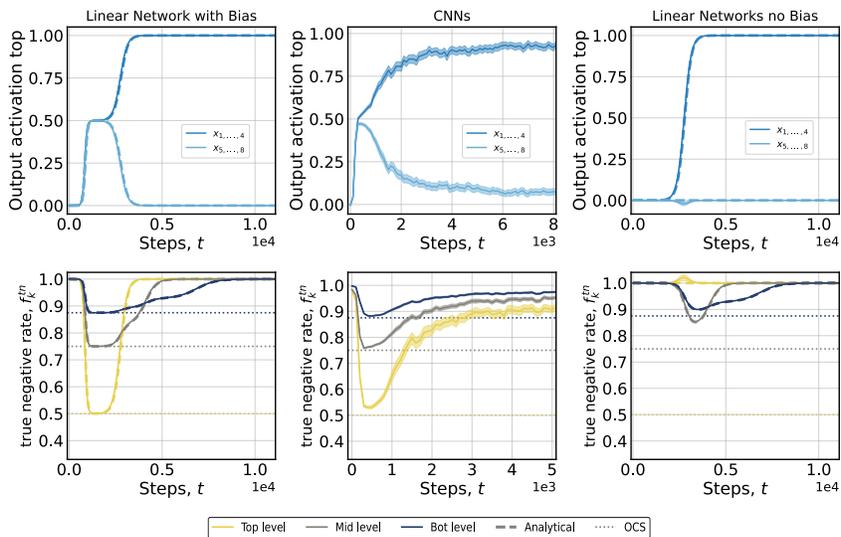


Figure 3: Early learning is driven to the OCS. *Top row*: Network predictions for a single output unit associated with the top level of the hierarchy in response to all inputs  $\mathbf{x}_i$ . We see clearly how CNNs and linear networks *with bias* initially change responses while not differentiating between different inputs before learning the correct input output mapping. *Bottom row*: True negative rates,  $f_k^{tn}$  for the three hierarchical levels as indicated by colors. For CNNs and linear networks *with bias* Performance approaches levels expected under the OCS (dotted lines).

the default one-hot label corresponding to each class  $i$  with the corresponding hierarchical, “three-hot” label  $\mathbf{y}_i$  seen in Fig. 1B. We use standard uniform Xavier initialization (Glorot and Bengio, 2010) and trained CNNs on an squared error loss. A full description of the CNN experiment and hyperparameter settings is deferred to Appendix A.8. We there also replicate our results with CIFAR-10 (Krizhevsky, 2009), non-hierarchical tasks, alternative loss functions, and CelebA (Liu et al., 2015) in Appendix A.9. We also show results for shallow networks in Appendix A.6.

To assess OCS learning we calculate true negative rates  $f_k^{tn}(\mathbf{y}, \hat{\mathbf{y}}) = \frac{(\mathbf{1}_N - \hat{\mathbf{y}}_k)^T (\mathbf{1}_N - \mathbf{y}_k)}{(\mathbf{1}_N - \mathbf{y}_k)^T (\mathbf{1}_N - \mathbf{y}_k)}$  for our task where the subscript  $k$  selects the vector slice corresponding to level  $k$  of the output hierarchy. We calculate the metric separately for the three hierarchical levels. Effectively, the metric describes how strongly model predictions  $\hat{\mathbf{y}}$  align with the desired outputs  $\mathbf{y}$  while focusing on zero entries only. The use of the metric is motivated by our desire to highlight how OCS learning is dependent on the distribution of labels in  $\mathbf{Y}$  and effectively measures wrong beliefs about the presence of target labels across the different levels of the hierarchy. Furthermore, the metric enables later comparisons to human learners (further details in Appendix A.4).

#### 4.1 EMPIRICAL EVIDENCE

We identify three separate empirical observations that support early learning of the OCS:

**Indifference.** Linear networks and CNNs initially change outputs while not differentiating between input examples. In Fig. 3 (top) we show the empirical and analytical activation of an output unit associated with the highest level of the hierarchy for all  $\mathbf{x}_i$ . Networks with and without bias terms learn to differentiate inputs correctly. However, networks with bias terms produce input-independent, non-zero outputs in early training as would be expected under the OCS.

**Performance.** Networks with bias terms show an initial tendency to over-select labels associated with the top level of the hierarchy as seen in the true negative rate  $f_k^{tn}(\mathbf{y}, \hat{\mathbf{y}})$  in Fig. 3 (bottom). Furthermore, linear networks and CNNs with bias terms almost exactly approach performance levels that would be provided by the OCS (dotted lines) for each of the three hierarchical levels. Linear network without bias terms do not produce this behaviour.

**OCS alignment.** The distance between outputs  $\hat{y}_i$  of linear and non-linear networks and  $y_{ocs}$  approaches zero in early training. Fig. 5 (top) shows how the  $L_1$  distance of sample-averaged network outputs and the OCS approaches zero before later converging to the desired network function.

## 4.2 THEORETICAL EXPLANATION

In this section, we extend the linear network formalism to understand the mechanism behind early learning of the OCS. We first show how bias terms in the input layer are directly related to the OCS. Afterwards, we prove that the OCS is necessarily learned first in these settings.

**The OCS is linked to shared properties.** Having established the applicability of the linear network theory in Section 3 we now seek to understand how the early bias towards the OCS emerges. To this end, notice how bias terms can be written in terms of the constant eigenmode  $\mathbf{1}_N$ :

**Proposition 2** (The OCS is linked to shared properties). *If  $\mathbf{1}_N$  is an eigenvector to the similarity matrix  $\mathbf{X}^T \mathbf{X} \in \mathbb{R}^{N \times N}$ , then the sample-average  $\bar{\mathbf{x}} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i$  will be an eigenvector to the correlation matrix  $\mathbf{X} \mathbf{X}^T \in \mathbb{R}^{N_{in} \times N_{in}}$  with identical eigenvalue  $\lambda$ . An analogous statement applies for  $\mathbf{Y}^T \mathbf{Y}$  and  $\mathbf{Y} \mathbf{Y}^T$ . The converse does not hold true in general.*

We prove this statement in Appendix A.5.5. Importantly, it establishes a connection between the feature and sample dimensions of  $\mathbf{X}$  and  $\mathbf{Y}$ . If  $\mathbf{1}_N$  is an eigenvector to  $\mathbf{X}^T \mathbf{X}$  and  $\mathbf{Y}^T \mathbf{Y}$  already, it implies that the addition of a bias term will directly add to its eigenvalue,  $s_{ocs}^2 \rightarrow s_{ocs}^2 + 1$ , even if it is initialized at zero. We show in Appendix A.5.6 that these assumptions on  $\mathbf{X}$  and  $\mathbf{Y}$  hold strictly for our hierarchical task design, and more generally relate to symmetry in the data (Appendix A.5.6). We discuss in Section 6 how this property extends to natural datasets where exact symmetry is absent.

Crucially, it now follows from Proposition 2 that the time-dependent network correlation  $\hat{\Sigma}^{yx}(t) = \mathbf{U} \mathbf{A}(t) \mathbf{V}^T$  in Eq. (2) will contain a strongly amplified OCS mode  $a_{ocs}(t) \mathbf{u}_{ocs} \mathbf{v}_{ocs}^T = a_{ocs}(t) \bar{\mathbf{y}} \bar{\mathbf{x}}^T$  by virtue of the modified singular value  $\sqrt{s_{ocs}^2 + 1}$  entering Eq. (2) and thereby the network function. Consequently, learning dynamics will be driven by the outer product of average input and output data. Moreover, this implies that given some input  $\mathbf{x}_i$  to Eq. (2), the network’s OCS mode contributes

$$\hat{y}_{ocs}(\mathbf{x}_i) = a_{ocs}(t) \mathbf{u}_{ocs} \mathbf{v}_{ocs}^T \mathbf{x}_i = a_{ocs}(t) \bar{\mathbf{y}} \bar{\mathbf{x}}^T \mathbf{x}_i \propto \bar{\mathbf{y}}. \quad (5)$$

The OCS mode in the time-dependent network function will hence necessarily drive responses towards average output statistics. Note that Eq. (5) also highlights that the more an input example is aligned to average inputs, the more the network’s responses will reflect average outputs. In particular, this makes the expected output  $\mathbb{E}_{\mathbf{x}}[\hat{y}(\mathbf{x})] \propto \bar{\mathbf{y}}$ . Throughout learning, the evolution of  $a_{ocs}(t)$  and scale-dependent alignment of  $\mathbf{x}_i$  and  $\bar{\mathbf{x}}$  will determine the network’s reliance on the OCS mode.

**Early learning is biased by the OCS mode.** We established that network responses are driven by average output statistics  $\bar{\mathbf{x}}$  and  $\bar{\mathbf{y}}$ , but why are *early* dynamics in particular influenced by the OCS? The learning speed of the SVD modes in the time-dependent network function are controlled by the magnitude of singular values  $s_\alpha$  as seen in Eq. (3).

**Theorem 1** (Early learning is biased by the OCS mode). *If  $\mathbf{1}_N$  is a joint non-degenerate eigenvector to positive input and output similarity matrices  $\mathbf{X}^T \mathbf{X}$  and  $\mathbf{Y}^T \mathbf{Y}$ , the OCS mode  $s_{ocs} \bar{\mathbf{y}} \bar{\mathbf{x}}^T$  will have leading spectral weight  $s_0 \equiv s_{ocs}$  in the SVD of the input-output correlation matrix  $\Sigma^{yx}$ .*

We prove this statement with help of the Perron-Frobenius theorem (Perron, 1907) in Appendix A.5.7. Consequently, the optimal constant mode is learned at a faster rate than remaining SVD components and transiently dominates the early network function. Notably, this applies to our task data  $\mathbf{Y}^T \mathbf{Y}$  (see Appendix A.5.6) and leads to characteristic learning signatures observed in Fig. 3.

Theorem 1 hinges on the constant eigenvector  $\mathbf{1}_N$  being present in the data. We later provide empirical (Fig. 6 and Appendix A.9) and theoretical (Appendix A.5.6) arguments that this assumption is approximately fulfilled in a variety of cases.

To recapitulate this section: We first rephrased a learnable bias term in the architecture as a shared feature in the input data. We then found that the associated singular value in Eq. (2) drives the learned network function towards the OCS (Eq. (5)). Finally, we proved that the bias affects *early* learning. In Appendix A.5.3, we summarize these results through the neural tangent kernel. Overall, these results demonstrate how architectural bias terms induce early OCS learning.

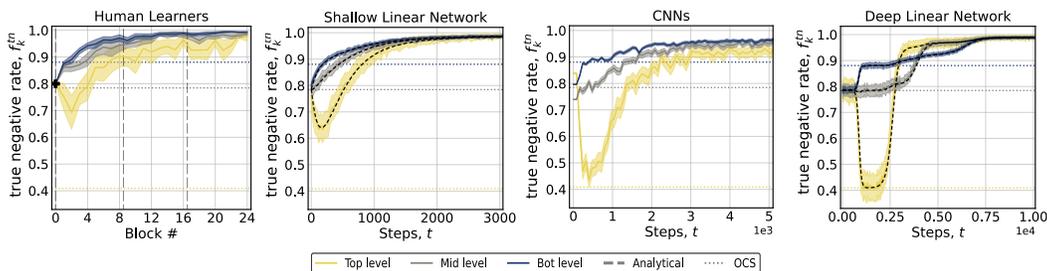


Figure 4: Early response bias towards the OCS across learners in the hierarchical learning task. True negative rates,  $f_k^{tn}$  for the three hierarchical levels as indicated by colors for biological and artificial learners (with bias terms). Dotted lines represent performances expected under the OCS. Observe that all learners show a transient bias towards the OCS. Dashed vertical gray lines indicates breaks between days for human learners.

## 5 SIMILARITIES BETWEEN LINEAR NETWORKS, HUMANS, AND COMPLEX MODELS

In this section, we demonstrate how human learners, linear networks, and non-linear architectures show strong similarities in their early learning on the hierarchical task displayed in Fig. 1.

**Setup.** The hierarchical learning task has previously been used extensively in the study of semantic cognition (Rogers and McClelland, 2004) and requires learners to develop a hierarchical one-to-many mapping as seen in Fig. 1B. We adapted the task for human learners while maintaining the underlying structure: Input stimuli were represented as different classes of planets and output labels were represented as a set of plant images (see Fig. 1E and Fig. 7). We also trained CNNs as in Section 4. Importantly, the hierarchical structure results in a non-uniform distribution of labels with average labels equal to  $y_{ocs}$ . Human learners received supervised training over three days. A full description of the experimental paradigm is given Appendix A.2. We then compute true negative rate  $f_k^{tn}(y, \hat{y})$  as in Section 4 while splitting performance across the hierarchical levels as before.

Neural networks produced continuous outputs in  $\mathbb{R}^{N_{out}}$  while humans responded via discrete button clicks in  $\{0, 1\}^{N_{out}}$ . As we do not have access to "human logits" before response execution we discretized network responses to enable comparison. We treat network responses in  $\hat{y}$  as logits from which we then sampled responses in  $\{0, 1\}^{N_{out}}$ . Full procedure details are given in Appendix A.3.

**Results.** The key results of our experiments are presented in Fig. 4. Intriguingly, we find that human learners, linear networks, and CNNs all display characteristic early response biases. Note that chance true negative rate is equal between all three levels of the hierarchy. Biological as well as artificial learners display an initial "drop" in true negative rate at the top level of the output hierarchy. The result indicates a general lack of specificity and an overly liberal response criterion for output labels on the top level of the hierarchy. To appreciate the significance of this result it is important to understand that the task can be learned without the development of these early response biases: In particular, linear networks without bias terms do not show this behaviour (see Appendix A.7). Surprisingly, the human response signature demonstrates that these learners, just as artificial networks, display an early bias towards the OCS. We conjecture that early learning of the OCS might be a general phenomenon that emerges during error-corrective training. We replicate the human result with a second cohort of learners in Appendix A.2. Notable is also the difference between shallow and deep linear networks. Response biases seem more transient in shallow networks and appear to more closely mirror human learners. However, quantitative comparisons are challenging due to inherently differing learning timescales.

## 6 GENERIC INPUT CORRELATIONS CAN EQUIVALENTLY DRIVE OCS LEARNING

We have established how the earliest phase of learning in linear networks is driven by the OCS. Crucially, in linear networks OCS learning hinges on bias terms in the network architecture. However,

432  
433  
434  
435  
436  
437  
438  
439  
440  
441  
442  
443  
444  
445  
446  
447  
448  
449  
450  
451  
452  
453  
454  
455  
456  
457  
458  
459  
460  
461  
462  
463  
464  
465  
466  
467  
468  
469  
470  
471  
472  
473  
474  
475  
476  
477  
478  
479  
480  
481  
482  
483  
484  
485

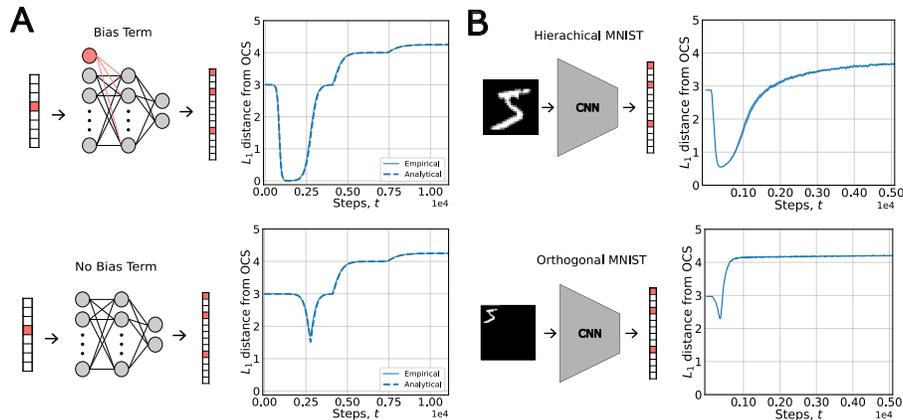


Figure 5:  $L_1$  Distance from the OCS in Linear networks and CNNs. **A** Linear networks with (top) and without bias terms (bottom) trained on the hierarchical task **B** CNNs without bias terms trained on variants of the hierarchical task. *Top*: Normal inputs. *Bottom*: "orthogonal" image inputs which remove all between-class input correlations from the input data. Note, how CNNs do not learn the OCS in the absence of input correlations and bias terms.

in non-linear architectures, such as CNNs, the network is driven towards the OCS even in the absence of bias terms (Fig. 5B, Top). The appearance of the data term  $\mathbf{X}^T \mathbf{X}$  in Proposition 3 suggests an equivalent effect that is induced by the data itself.

**Corollary 1** (Input correlations induce early OCS). *If  $\mathbf{1}_N$  is an eigenvector of the data similarity matrix  $\mathbf{X}^T \mathbf{X}$  with non-degenerate eigenvalue  $s_0$ , then the OCS response during early learning will be driven according to its magnitude.*

This statement follows directly from the joint diagonalisation of Eq. (1) and subsequent projection onto the OCS  $\mathbf{y}_{OCS}$ . We show a solvable case of OCS learning in linear networks under input correlations and in the absence of bias terms in Appendix A.11. We furthermore hypothesize that neural networks will be driven towards the OCS if training data contains more generic input correlations where  $\mathbf{1}_N$  is not an exact eigenvector.

**Setup.** We trained CNNs on the hierarchical task in Section 4. Inputs were given by eight randomly sampled classes of MNIST (Fig. 5B, Top). To isolate the effect of input correlations we created a second dataset where randomly sampled classes of MNIST were copied on orthogonal subspaces of a larger image (Fig. 5B, bottom). Importantly, this procedure removes all between-class correlations.

**Results.** The main result of our experiment is displayed in Fig. 5. CNNs which learn from standard MNIST images are strongly driven towards the OCS. In contrast, early dynamics for the "orthogonal" MNIST do not display this tendency. Strikingly, the early dynamics with standard MNIST classes are highly similar to those observed in linear network with bias terms, while the dynamics for the latter task resemble those seen in the linear network without this feature. To verify that generic input correlations are indeed causing these differences we explore the eigenspectrum of the data correlation matrices. We sample 100 images from all 10 classes and compute a correlation matrix  $\mathbf{X}^T \mathbf{X}$  from flattened images. First, we find that the eigenspectrum for standard MNIST images is dominated by a single eigenvector (Fig. 6, top-left). In contrast, the eigenspectrum of the orthogonal MNIST task does not display this property (Fig. 6, top-center). Further, recall that input bias terms lead to a non-degenerate constant eigenvector  $\mathbf{1}_N$  in the input correlation matrix (Section 4). Similarly, we find that the first eigenvector  $\mathbf{v}_1$  of  $\mathbf{X}^T \mathbf{X}$  is indeed highly aligned to  $\mathbf{1}_N$  (Fig. 6, right), whereas this is not the case in the orthogonal MNIST. We additionally show similar results for CIFAR-10 and CelebA. Theoretical considerations suggest that these correlations originate from an approximate symmetry in the data (see Appendix A.5.6).

Overall, we here demonstrated that early learning of the OCS can be driven by properties of the architecture (bias terms) or data (input correlations). Our results also highlight that input correlations are a common feature of standard image datasets: Early learning of the OCS might be a common

occurrence when learning from such data. To see a practical implication of these results we briefly discuss fairness implications of OCS learning in Appendix A.10.

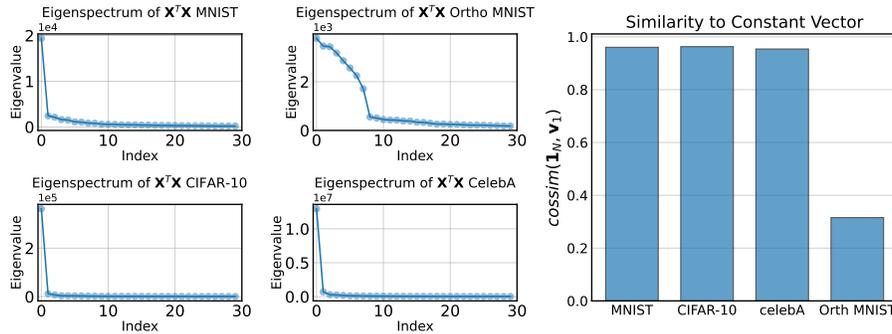


Figure 6: Dataset eigenspectra and constancy of first eigenvector for different image datasets. *Left:* Eigenspectra of  $X^T X$  for different datasets. *Right:* Alignment of first eigenvector in  $X^T X$  with the constant vector  $\mathbf{1}_N$ .

## 7 DISCUSSION

In this work, we found that the inclusion of bias terms in linear networks shifts early learning towards the OCS, even when initialized at zero. We also highlight how OCS learning can equivalently be driven by input correlations. We demonstrated that early, input-independent simplicity biases occur in practice, affecting both non-linear networks and human learners. Our contribution complements prior work on simplicity biases by highlighting factors that drive networks in the *earliest* stages of learning; connecting input, output, and architecture. Overall, our findings highlight how simple linear networks can serve as useful tools to investigate simplicity biases in significantly more complex systems.

**Relevance.** We see promising applications for early OCS learning in the cognitive and behavioural sciences. OCS-like response biases have been noted previously (Herrnstein, 1961; Estes, 1964). However, we believe that a normative theory for these effects is still incomplete. Our theory identifies possible properties of the biological wetware or natural stimuli that may give rise to such biases.

While we do not study generalisation ourselves, we believe that OCS learning is practically important to understand *how* neural networks generalize or fail to generalize. Kang et al. (2024) has highlighted that networks will revert to OCS in a variety of generalisation settings. We demonstrate that the OCS component in the network function is acquired *early*, and is *retained* throughout training (effective singular values in Fig. 2 stay constant in late training). We believe that this retention of the OCS mode enables reversion.

OCS learning is also relevant when learning under class imbalance, a common problem in machine learning where datasets are frequently naturally imbalanced (Feldman, 2020; Van Horn and Perona, 2017), leading to a failure to learn information about minority classes (Ye et al., 2021). In Appendix A.10 we show an exactly solvable case of OCS learning in such settings and highlight how OCS learning can negatively impact performance for minority classes.

**Limitations and future work.** Our work is restricted to qualitative comparisons between linear networks and non-linear systems and our work only gives suggestive evidence of factors which drive early OCS learning in non-linear systems. We chose linear networks to allow for a rigorous description of the dynamics of learning. Methods from mean-field theory may provide a precise tool to analyze a wider range of systems directly. Second, the ambiguity between architecture and data in driving the OCS does not allow us to determine the underlying mechanism in human learners. Future studies might address this limitation by manipulating correlations in stimuli or by recording of neural data.

**Reproducibility statement.** We provide the code to produce our simulation results in the supplementary material to this submission.

540  
541  
542  
543  
544  
545  
546  
547  
548  
549  
550  
551  
552  
553  
554  
555  
556  
557  
558  
559  
560  
561  
562  
563  
564  
565  
566  
567  
568  
569  
570  
571  
572  
573  
574  
575  
576  
577  
578  
579  
580  
581  
582  
583  
584  
585  
586  
587  
588  
589  
590  
591  
592  
593

---

## REFERENCES

- Dimitris Kalimeris, Gal Kaplun, Preetum Nakkiran, Benjamin Edelman, Tristan Yang, Boaz Barak, and Haofeng Zhang. SGD on Neural Networks Learns Functions of Increasing Complexity. In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- Nasim Rahaman, Aristide Baratin, Devansh Arpit, Felix Draxler, Min Lin, Fred Hamprecht, Yoshua Bengio, and Aaron Courville. On the Spectral Bias of Neural Networks. In *Proceedings of the 36th International Conference on Machine Learning*, pages 5301–5310. PMLR, May 2019. URL <https://proceedings.mlr.press/v97/rahaman19a.html>. ISSN: 2640-3498.
- Satwik Bhattamishra, Arkil Patel, Varun Kanade, and Phil Blunsom. Simplicity Bias in Transformers and their Ability to Learn Sparse Boolean Functions, July 2023. URL <http://arxiv.org/abs/2211.12316>. arXiv:2211.12316 [cs].
- Guillermo Valle-Pérez, Chico Q. Camargo, and Ard A. Louis. Deep learning generalizes because the parameter-function map is biased towards simple functions, April 2019. URL <http://arxiv.org/abs/1805.08522>. arXiv:1805.08522 [cs, stat].
- Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning (still) requires rethinking generalization. *Communications of the ACM*, 64(3): 107–115, March 2021. ISSN 0001-0782, 1557-7317. doi: 10.1145/3446776. URL <https://dl.acm.org/doi/10.1145/3446776>.
- Maria Refinetti, Alessandro Ingrosso, and Sebastian Goldt. Neural networks trained with SGD learn distributions of increasing complexity. In *Proceedings of the 40th International Conference on Machine Learning*, pages 28843–28863. PMLR, July 2023. URL <https://proceedings.mlr.press/v202/refinetti23a.html>. ISSN: 2640-3498.
- Nora Belrose, Quintin Pope, Lucia Quirke, Alex Mallen, and Xiaoli Fern. Neural Networks Learn Statistics of Increasing Complexity, February 2024. URL <http://arxiv.org/abs/2402.04362>. arXiv:2402.04362 [cs].
- Lukas Braun, Clémentine Dominé, James Fitzgerald, and Andrew Saxe. Exact learning dynamics of deep linear networks with prior knowledge. *Advances in Neural Information Processing Systems*, 35:6615–6629, December 2022. URL [https://proceedings.neurips.cc/paper\\_files/paper/2022/hash/2b3bb2c95195130977a51b3bb251c40a-Abstract-Conference.html](https://proceedings.neurips.cc/paper_files/paper/2022/hash/2b3bb2c95195130977a51b3bb251c40a-Abstract-Conference.html).
- Andrew M. Saxe, James L. McClelland, and Surya Ganguli. Exact solutions to the nonlinear dynamics of learning in deep linear neural networks, February 2014. URL <http://arxiv.org/abs/1312.6120>. arXiv:1312.6120 [cond-mat, q-bio, stat].
- Andrew M. Saxe, James L. McClelland, and Surya Ganguli. A mathematical theory of semantic development in deep neural networks. *Proceedings of the National Academy of Sciences*, 116(23): 11537–11546, June 2019. doi: 10.1073/pnas.1820226116. URL <https://www.pnas.org/doi/full/10.1073/pnas.1820226116>. Publisher: Proceedings of the National Academy of Sciences.
- Timothy T. Rogers and James L. McClelland. *Semantic Cognition: A Parallel Distributed Processing Approach*. MIT Press, 2004. ISBN 978-0-262-18239-3. Google-Books-ID: AmB33Uz2MVAC.
- David E. Rumelhart, James L. McClelland, and PDP Research Group. *Parallel distributed processing, volume 1: Explorations in the microstructure of cognition: Foundations*. The MIT press, 1986. URL <https://scholar.google.com/scholar?cluster=13839636846206420541&hl=en&oi=scholarrr>.
- Katie Kang, Amrith Setlur, Claire Tomlin, and Sergey Levine. Deep Neural Networks Tend To Extrapolate Predictably, March 2024. URL <http://arxiv.org/abs/2310.00873>. arXiv:2310.00873 [cs].

- 
- 594 R. J. Herrnstein. Relative and absolute strength of response as a function of frequency of rein-  
595 forcement,. *Journal of the Experimental Analysis of Behavior*, 4(3):267–272, July 1961. ISSN  
596 0022-5002. doi: 10.1901/jeab.1961.4-267. URL [https://www.ncbi.nlm.nih.gov/pmc/  
597 articles/PMC1404074/](https://www.ncbi.nlm.nih.gov/pmc/articles/PMC1404074/).
- 598 William K. Estes. Probability Learning and Sequence learning. In Arthur W. Melton, editor,  
599 *Categories of Human Learning*, pages 89–128. Academic Press, January 1964. ISBN 978-1-4832-  
600 3145-7. doi: 10.1016/B978-1-4832-3145-7.50010-8. URL [https://www.sciencedirect.  
601 com/science/article/pii/B9781483231457500108](https://www.sciencedirect.com/science/article/pii/B9781483231457500108).
- 602 W. K. Estes and J. H. Straughan. Analysis of a verbal conditioning situation in terms of statistical  
603 learning theory. *Journal of Experimental Psychology*, 47(4):225–234, 1954. ISSN 0022-1015. doi:  
604 10.1037/h0060989. URL <https://doi.apa.org/doi/10.1037/h0060989>.
- 605 Pete R. Jones, David R. Moore, Daniel E. Shub, and Sygal Amitay. The role of response bias in  
606 perceptual learning. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, 41  
607 (5):1456–1470, September 2015. ISSN 1939-1285, 0278-7393. doi: 10.1037/xlm0000111. URL  
608 <https://doi.apa.org/doi/10.1037/xlm0000111>.
- 609 Joshua I. Gold, Chi-Tat Law, Patrick Connolly, and Sharath Bennur. The Relative Influences of  
610 Priors and Sensory Evidence on an Oculomotor Decision Variable During Perceptual Learning.  
611 *Journal of Neurophysiology*, 100(5):2653–2668, November 2008. ISSN 0022-3077, 1522-1598.  
612 doi: 10.1152/jn.90629.2008. URL [https://www.physiology.org/doi/10.1152/jn.  
613 90629.2008](https://www.physiology.org/doi/10.1152/jn.90629.2008).
- 614 William S. Verplanck, George H. Collier, and John W. Cotton. Nonindependence of successive  
615 responses in measurements of the visual threshold. *Journal of Experimental Psychology*, 44(4):  
616 273–282, 1952. ISSN 0022-1015. doi: 10.1037/h0054948. URL [https://doi.apa.org/  
617 doi/10.1037/h0054948](https://doi.apa.org/doi/10.1037/h0054948).
- 618 James R. Hawker. The influence of training procedure and other task variables in paired-associate  
619 learning. *Journal of Verbal Learning and Verbal Behavior*, 3(1):70–76, February 1964. ISSN  
620 0022-5371. doi: 10.1016/S0022-5371(64)80060-8. URL [https://www.sciencedirect.  
621 com/science/article/pii/S0022537164800608](https://www.sciencedirect.com/science/article/pii/S0022537164800608).
- 622 Gordon H. Bower. An association model for response and training variables in paired-associate  
623 learning. *Psychological Review*, 69(1):34–53, January 1962. ISSN 1939-1471, 0033-295X. doi:  
624 10.1037/h0039023. URL <https://doi.apa.org/doi/10.1037/h0039023>.
- 625 Jacob Feldman. Minimization of Boolean complexity in human concept learning. *Nature*, 407  
626 (6804):630–633, October 2000. ISSN 0028-0836, 1476-4687. doi: 10.1038/35036586. URL  
627 <https://www.nature.com/articles/35036586>.
- 628 Noah D. Goodman, Joshua B. Tenenbaum, Jacob Feldman, and Thomas L. Griffiths.  
629 A Rational Analysis of Rule-Based Concept Learning. *Cognitive Science*, 32(1):108–  
630 154, 2008. ISSN 1551-6709. doi: 10.1080/03640210701802071. URL [https://  
631 onlinelibrary.wiley.com/doi/abs/10.1080/03640210701802071](https://onlinelibrary.wiley.com/doi/abs/10.1080/03640210701802071). \_eprint:  
632 <https://onlinelibrary.wiley.com/doi/pdf/10.1080/03640210701802071>.
- 633 Nick Chater. Reconciling Simplicity and Likelihood Principles in Perceptual Organization. *Psycho-*  
634 *logical review*, 103:566–81, July 1996. doi: 10.1037/0033-295X.103.3.566.
- 635 Tania Lombrozo. Simplicity and probability in causal explanation. *Cognitive Psychology*, 55(3):  
636 232–257, November 2007. ISSN 0010-0285. doi: 10.1016/j.cogpsych.2006.09.006. URL [https://  
637 www.sciencedirect.com/science/article/pii/S0010028506000739](https://www.sciencedirect.com/science/article/pii/S0010028506000739).
- 638 Jacob Feldman. The Simplicity Principle in Human Concept Learning. *Current Directions in Psycho-*  
639 *logical Science*, 12(6):227–232, December 2003. ISSN 0963-7214. doi: 10.1046/j.0963-7214.2003.  
640 01267.x. URL <https://doi.org/10.1046/j.0963-7214.2003.01267.x>. Pub-  
641 lisher: SAGE Publications Inc.
- 642 Kenji Fukumizu. Effect Of Batch Learning In Multilayer Neural Networks. June 1998.

---

648 Pierre Baldi and Kurt Hornik. Neural networks and principal component analysis: Learning from  
649 examples without local minima. *Neural Networks*, 2(1):53–58, January 1989. ISSN 0893-6080. doi:  
650 10.1016/0893-6080(89)90014-2. URL [https://www.sciencedirect.com/science/  
651 article/pii/0893608089900142](https://www.sciencedirect.com/science/article/pii/0893608089900142).

652 Andrew K. Lampinen and Surya Ganguli. An analytic theory of generalization dynamics and  
653 transfer learning in deep linear networks, January 2019. URL [http://arxiv.org/abs/  
654 1809.10374](http://arxiv.org/abs/1809.10374). arXiv:1809.10374 [cs, stat].

655 Samuel Liebana Garcia, Aeron Laffere, Chiara Toschi, Louisa Schilling, Jacek Podlaski, Matthias  
656 Fritsche, Peter Zátka-Haas, Yulong Li, Rafal Bogacz, Andrew Saxe, and Armin Lak. Striatal  
657 dopamine reflects individual long-term learning trajectories, December 2023. URL [http://  
658 biorxiv.org/lookup/doi/10.1101/2023.12.14.571653](http://biorxiv.org/lookup/doi/10.1101/2023.12.14.571653).

659 Sygal Amitay, Yu-Xuan Zhang, Pete R. Jones, and David R. Moore. Perceptual learning: Top  
660 to bottom. *Vision Research*, 99:69–77, June 2014. ISSN 0042-6989. doi: 10.1016/j.visres.  
661 2013.11.006. URL [https://www.sciencedirect.com/science/article/pii/  
662 S0042698913002800](https://www.sciencedirect.com/science/article/pii/S0042698913002800).

663 Anne E Urai, Jan Willem de Gee, Konstantinos Tsetsos, and Tobias H Donner. Choice history  
664 biases subsequent evidence accumulation. *eLife*, 8:e46331, July 2019. ISSN 2050-084X. doi:  
665 10.7554/eLife.46331. URL <https://doi.org/10.7554/eLife.46331>. Publisher: eLife  
666 Sciences Publications, Ltd.

667 Yunshu Fan, Takahiro Doi, Joshua I. Gold, and Long Ding. Neural Representations of Post-Decision  
668 Accuracy and Reward Expectation in the Caudate Nucleus and Frontal Eye Field. *The Journal  
669 of Neuroscience*, 44(2):e0902232023, January 2024. ISSN 0270-6474, 1529-2401. doi: 10.  
670 1523/JNEUROSCI.0902-23.2023. URL [https://www.jneurosci.org/lookup/doi/  
671 10.1523/JNEUROSCI.0902-23.2023](https://www.jneurosci.org/lookup/doi/10.1523/JNEUROSCI.0902-23.2023).

672 Leo P. Sugrue, Greg S. Corrado, and William T. Newsome. Matching Behavior and the Representation  
673 of Value in the Parietal Cortex. *Science*, 304(5678):1782–1787, June 2004. ISSN 0036-8075,  
674 1095-9203. doi: 10.1126/science.1094765. URL [https://www.science.org/doi/10.  
675 1126/science.1094765](https://www.science.org/doi/10.1126/science.1094765).

676 Gilles Dutilh, Don van Ravenzwaaij, Sander Nieuwenhuis, Han L. J. van der Maas, Birte U.  
677 Forstmann, and Eric-Jan Wagenmakers. How to measure post-error slowing: A confound and a  
678 simple solution. *Journal of Mathematical Psychology*, 56(3):208–216, June 2012. ISSN 0022-2496.  
679 doi: 10.1016/j.jmp.2012.04.001. URL [https://www.sciencedirect.com/science/  
680 article/pii/S0022249612000454](https://www.sciencedirect.com/science/article/pii/S0022249612000454).

681 Patrick Rabbitt and Bryan Rodgers. What does a Man do after he Makes an Error? An  
682 Analysis of Response Programming. *Quarterly Journal of Experimental Psychology*, 29(4):  
683 727–743, November 1977. ISSN 0033-555X. doi: 10.1080/14640747708400645. URL  
684 <http://journals.sagepub.com/doi/10.1080/14640747708400645>.

685 Blake Bordelon, Abdulkadir Canatar, and Cengiz Pehlevan. Spectrum dependent learning curves  
686 in kernel regression and wide neural networks. *ArXiv e-prints*, 2020. URL [https://arxiv.  
687 org/abs/2002.02561](https://arxiv.org/abs/2002.02561). tex.eprint: 2002.02561.

688 Song Mei, Theodor Misiakiewicz, and Andrea Montanari. Generalization error of random feature and  
689 kernel methods: hypercontractivity and kernel matrix concentration. *Applied and Computational  
690 Harmonic Analysis*, 59:3–84, 2022. Publisher: Elsevier tex.creationdate: 2022-07-20T21:54:34  
691 tex.modificationdate: 2022-07-20T21:54:42.

692 Chris Mingard, Henry Rees, Guillermo Valle-Pérez, and Ard A. Louis. Do deep neural networks  
693 have an inbuilt Occam’s razor?, April 2023. URL <http://arxiv.org/abs/2304.06670>.  
694 arXiv:2304.06670 [cs, stat].

695 Li Deng. The MNIST Database of Handwritten Digit Images for Machine Learning Research [Best of  
696 the Web]. *IEEE Signal Processing Magazine*, 29(6):141–142, November 2012. ISSN 1053-5888.  
697 doi: 10.1109/MSP.2012.2211477. URL [http://ieeexplore.ieee.org/document/  
698 6296535/](http://ieeexplore.ieee.org/document/6296535/).

- 
- 702 Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural  
703 networks. In *Proceedings of the Thirteenth International Conference on Artificial Intelligence*  
704 *and Statistics*, pages 249–256. JMLR Workshop and Conference Proceedings, March 2010. URL  
705 <https://proceedings.mlr.press/v9/glorot10a.html>. ISSN: 1938-7228.
- 706 Alex Krizhevsky. Learning Multiple Layers of Features from Tiny Images. <https://www.cs.toronto.edu/kriz/learning-features-2009-TR.pdf>, 2009.
- 707  
708
- 709 Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep Learning Face Attributes in the Wild,  
710 September 2015. URL <http://arxiv.org/abs/1411.7766>. arXiv:1411.7766 [cs].
- 711
- 712 Oskar Perron. Zur Theorie der Matrices. *Mathematische Annalen*, 64(2):248–263, June 1907. ISSN  
713 0025-5831, 1432-1807. doi: 10.1007/BF01449896. URL [http://link.springer.com/](http://link.springer.com/10.1007/BF01449896)  
714 [10.1007/BF01449896](http://link.springer.com/10.1007/BF01449896).
- 715 Vitaly Feldman. Does learning require memorization? a short tale about a long tail. In *Proceedings*  
716 *of the 52nd Annual ACM SIGACT Symposium on Theory of Computing*, pages 954–959, Chicago  
717 IL USA, June 2020. ACM. ISBN 978-1-4503-6979-4. doi: 10.1145/3357713.3384290. URL  
718 <https://dl.acm.org/doi/10.1145/3357713.3384290>.
- 719 Grant Van Horn and Pietro Perona. The Devil is in the Tails: Fine-grained Classification in the Wild,  
720 September 2017. URL <http://arxiv.org/abs/1709.01450>. arXiv:1709.01450 [cs].
- 721
- 722 Han-Jia Ye, De-Chuan Zhan, and Wei-Lun Chao. Procrustean Training for Imbalanced Deep Learning.  
723 In *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 92–102, Montreal,  
724 QC, Canada, October 2021. IEEE. ISBN 978-1-66542-812-5. doi: 10.1109/ICCV48922.2021.  
725 00016. URL <https://ieeexplore.ieee.org/document/9710650/>.
- 726 Arthur Jacot, Franck Gabriel, and Clement Hongler. Neural Tangent Kernel: Convergence and  
727 Generalization in Neural Networks. In *Advances in Neural Information Processing Systems*,  
728 volume 31. Curran Associates, Inc., 2018. URL [https://proceedings.neurips.cc/](https://proceedings.neurips.cc/paper/2018/hash/5a4be1fa34e62bb8a6ec6b91d2462f5a-Abstract.html)  
729 [paper/2018/hash/5a4be1fa34e62bb8a6ec6b91d2462f5a-Abstract.html](https://proceedings.neurips.cc/paper/2018/hash/5a4be1fa34e62bb8a6ec6b91d2462f5a-Abstract.html).
- 730 Daniel A Roberts, Sho Yaida, and Boris Hanin. The Principles of Deep Learning Theory. page 449.
- 731
- 732 Erich Hecke. Über orthogonal-invariante integralgleichungen. *Mathematische Annalen*, 78  
733 (1):398–404, 1917. Publisher: Springer-Verlag tex.creationdate: 2022-07-23T15:10:14  
734 tex.modificationdate: 2022-07-23T15:10:14.
- 735 Vincent Dutordoir, Nicolas Durrande, and James Hensman. Sparse Gaussian processes with spherical  
736 harmonic features. In *International Conference on Machine Learning*, pages 2793–2802. PMLR,  
737 2020. ISBN 2640-3498.
- 738
- 739 Ayşe Erzan and Aslı Tuncer. Explicit construction of the eigenvectors and eigenvalues of the graph  
740 Laplacian on the Cayley tree. *Linear Algebra and its Applications*, 586:111–129, February 2020.  
741 ISSN 0024-3795. doi: 10.1016/j.laa.2019.10.023. URL <https://www.sciencedirect.com/science/article/pii/S002437951930463X>.
- 742
- 743 Andries E. Brouwer and Willem H. Haemers. *Spectra of Graphs*. Springer Science & Business  
744 Media, December 2011. ISBN 978-1-4614-1939-6. Google-Books-ID: F98THwYgrXYC.
- 745
- 746 Haibo He and E.A. Garcia. Learning from Imbalanced Data. *IEEE Transactions on Knowledge and*  
747 *Data Engineering*, 21(9):1263–1284, September 2009. ISSN 1041-4347. doi: 10.1109/TKDE.  
748 2008.239. URL <http://ieeexplore.ieee.org/document/5128907/>.
- 749
- 750 Chen Huang, Yining Li, Chen Change Loy, and Xiaoou Tang. Learning Deep Representation for  
751 Imbalanced Classification. In *2016 IEEE Conference on Computer Vision and Pattern Recognition*  
752 *(CVPR)*, pages 5375–5384, Las Vegas, NV, USA, June 2016. IEEE. ISBN 978-1-4673-8851-  
753 1. doi: 10.1109/CVPR.2016.580. URL [https://ieeexplore.ieee.org/document/](https://ieeexplore.ieee.org/document/7780949/)  
754 [7780949/](https://ieeexplore.ieee.org/document/7780949/).
- 755 Emanuele Francizi, Marco Baity-Jesi, and Aurelien Lucchi. A Theoretical Analysis of the Learning  
Dynamics under Class Imbalance, June 2023. URL <http://arxiv.org/abs/2207.00391>.  
arXiv:2207.00391 [cs, stat].

---

756 Bin Liu, Konstantinos Blekas, and Grigorios Tsoumakas. Multi-Label Sampling based on Local Label  
757 Imbalance, May 2020. URL <http://arxiv.org/abs/2005.03240>. arXiv:2005.03240  
758 [cs].

759 Francisco Charte, Antonio J. Rivera, María J. del Jesus, and Francisco Herrera. Addressing imbalance  
760 in multilabel classification: Measures and random resampling algorithms. *Neurocomputing*, 163:  
761 3–16, September 2015. ISSN 0925-2312. doi: 10.1016/j.neucom.2014.08.091. URL [https://](https://www.sciencedirect.com/science/article/pii/S0925231215004269)  
762 [www.sciencedirect.com/science/article/pii/S0925231215004269](https://www.sciencedirect.com/science/article/pii/S0925231215004269).  
763

764 Khoi Pham, Kushal Kafle, Zhe Lin, Zhihong Ding, Scott Cohen, Quan Tran, and Abhinav Shrivastava.  
765 Learning to Predict Visual Attributes in the Wild. In *2021 IEEE/CVF Conference on Computer*  
766 *Vision and Pattern Recognition (CVPR)*, pages 13013–13023, Nashville, TN, USA, June 2021.  
767 IEEE. ISBN 978-1-66544-509-2. doi: 10.1109/CVPR46437.2021.01282. URL [https://](https://ieeexplore.ieee.org/document/9578060/)  
768 [ieeexplore.ieee.org/document/9578060/](https://ieeexplore.ieee.org/document/9578060/).

769 Weiwei Liu, Haobo Wang, Xiaobo Shen, and Ivor W. Tsang. The Emerging Trends of Multi-  
770 Label Learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(11):7955–  
771 7974, November 2022. ISSN 1939-3539. doi: 10.1109/TPAMI.2021.3119334. URL [https://](https://ieeexplore.ieee.org/abstract/document/9568738)  
772 [ieeexplore.ieee.org/abstract/document/9568738](https://ieeexplore.ieee.org/abstract/document/9568738). Conference Name: IEEE  
773 Transactions on Pattern Analysis and Machine Intelligence.

774 Yin Cui, Menglin Jia, Tsung-Yi Lin, Yang Song, and Serge Belongie. Class-Balanced Loss Based  
775 on Effective Number of Samples. pages 9268–9277, 2019. URL [https://openaccess.](https://openaccess.thecvf.com/content_CVPR_2019/html/Cui_Class-Balanced_Loss_Based_on_Effective_Number_of_Samples_CVPR_2019_paper.html)  
776 [thecvf.com/content\\_CVPR\\_2019/html/Cui\\_Class-Balanced\\_Loss\\_Based\\_](https://openaccess.thecvf.com/content_CVPR_2019/html/Cui_Class-Balanced_Loss_Based_on_Effective_Number_of_Samples_CVPR_2019_paper.html)  
777 [on\\_Effective\\_Number\\_of\\_Samples\\_CVPR\\_2019\\_paper.html](https://openaccess.thecvf.com/content_CVPR_2019/html/Cui_Class-Balanced_Loss_Based_on_Effective_Number_of_Samples_CVPR_2019_paper.html).  
778  
779  
780  
781  
782  
783  
784  
785  
786  
787  
788  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809

## A APPENDIX / SUPPLEMENTAL MATERIAL

### A.1 OVERVIEW

Our appendix has the following sections:

- In Appendix A.2, we describe the human experiment in more detail and show the results of a replication in a second cohort. We furthermore report statistical tests and describe ethical considerations.
- In Appendix A.3, we outline how we bring neural network and human responses displayed in Section 5 into a common space for direct comparison.
- In Appendix A.4, we outline how we compute the true negative rates,  $f^{tn}$  used in Section 4 and Section 5.
- In Appendix A.5, we provide additional theoretical derivations and remaining proofs to the statements in the main text.
- In Appendix A.6, we show OCS signatures in *shallow* networks with bias terms.
- In the short Appendix A.7, we show how linear networks without bias terms behave on the task in Section 5.
- In Appendix A.8, we describe hyperparameters, datasets, and further training details used for our CNN experiments.
- In Appendix A.9, we describe the results of additional experiments investigating early emergence of the OCS in non-linear models.
- In Appendix A.10 we show an additional solvable case of linear networks with bias terms under class imbalance.
- In Appendix A.11 we show OCS learning in linear networks with input correlations but in the absence of bias terms.
- In Appendix A.12 we discuss additional connections of our work to multi-label learning.

### A.2 HUMAN LEARNING EXPERIMENT

We directly translated the hierarchical task setup used by Saxe et al. (2019) into an experimental paradigm. Our design attempts to stay as close to the original task structure used for neural networks as possible. We designed the task as a mapping from 8 distinct input stimuli represented as planets to a set of 3 associated output stimuli represented as plants (see Fig. 7, left).

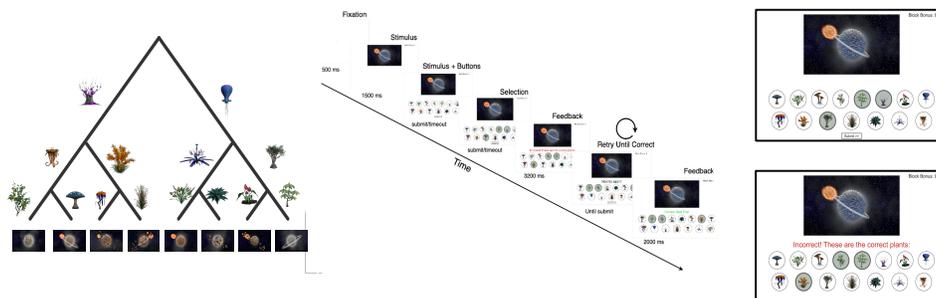


Figure 7: Human task design. *Left*: Hierarchical learning task, adapted for human participants. *Centre*: Trial structure as experienced by human participants. *Right*: Example screen during response period (top), Example screen during feedback period (bottom).

In the task, participants had to learn to associate which outputs properties are associated with each input. Unbeknownst to the participants we imposed a hierarchical structure on output targets (Fig. 7, left). In the structure some output labels are associated with more than one input. As a control for analyses we also included an additional control input-output pair (similarly represented by a planet

864  
865  
866  
867  
868  
869  
870  
871  
872  
873  
874  
875  
876  
877  
878  
879  
880  
881  
882  
883  
884  
885  
886  
887  
888  
889  
890  
891  
892  
893  
894  
895  
896  
897  
898  
899  
900  
901  
902  
903  
904  
905  
906  
907  
908  
909  
910  
911  
912  
913  
914  
915  
916  
917

and a plant; not shown here and excluded from current analysis). We recruited a cohort of 10 subjects that were trained over the course of three days with one daily session. The cohort was recruited as part of a larger neuroimaging experiment but our analysis presented here is exclusive to behavioural results. We further replicated our results in a second cohort of 46 human subjects recruited via the online platform Prolific (prolific.com). Results of the replication of the study can be seen in Fig. 8.

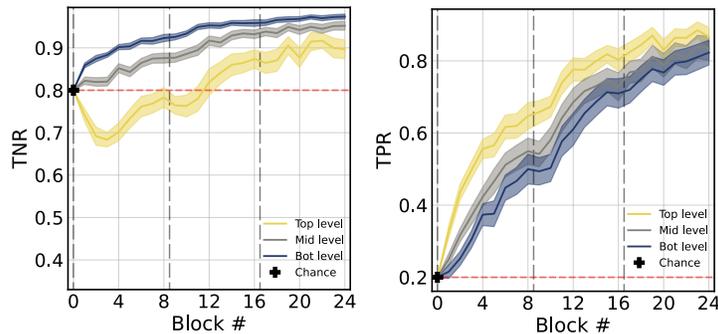


Figure 8: The human replication cohort. While learning is slower, the qualitative pattern indicating reliance on the OCS is replicated. *Left*: TNR rate. *Right*: TPR rate.

Each day of at home training consisted of 8 blocks of training with 22 trials each (160 standard trials and 16 control trials) which lasted about one hour. The trial structure during training is shown in Fig. 7, centre. During training trials, subjects were shown the stimulus on screen and were required to press three buttons, presented below the planet image (Fig. 7, right). The subjects received fully informative feedback on each trial and were forced to repeat the trial in the case of incorrect selection until the correct properties were selected. The location of buttons was shuffled on screen for each trial and for each forced repetition. For each button clicked correctly on their first attempt participants received a bonus point. We displayed a block-wise bonus in the corner of the screen throughout the task. Participants were payed slightly above local minimum wage as a baseline and received a substantial performance dependent bonus (on average about one-third of the baseline pay). We include a screenshot of the initial instructions in Fig. 9. Beyond this initial instruction screen participants received more nuanced instructions about clicking of buttons and feedback in the beginning of the task.

**Statistical tests.** While our focus is on qualitative patterns in human behaviour, we compute statistical tests on the true negative rates for human results seen in the main text (Fig. 4). We averaged all blocks in a given day and performed a two-way repeated measures ANOVA to assess the effect of day and hierarchy level on true negative rates. The two-way repeated measures ANOVA revealed significant main effects of day  $F(2, 18) = 57.22, p < .0001, \eta^2 = .25$  and level  $F(2, 18) = 6.25, p = .033, \eta^2 = .18$ . Beyond this we also found a significant interaction of day and level  $F(4, 36) = 9.795, p = .0056, \eta^2 = .042$ . A Mauchly test indicated that the assumption of sphericity had been violated for level  $\chi^2 = .03, p < .5$  and the interaction term  $\chi^2 = .006, p < .5$ . Significance values are reported with Greenhouse-Geisser correction. The results confirm that performance between levels are significantly different depending on day and hierarchical level.

**Ethical considerations.** Human participants performed a simple, computerised learning task without the collection of personal identifiable information or substantial deception. Human data collection was handled strictly in line with institutional guidelines and under institutional review board approval. We obtained informed consent for each participant before commencing the study. We highlighted that participants could withdraw at any time without penalty or loss of compensation by simply exiting full-screen or informing the experimenter. We provided contact emails in the case of concern or questions. Data was handled in a strictly anonymised format and stored on password secured devices. Participants were payed above minimum wage for their country of origin.

918  
919  
920  
921  
922  
923  
924  
925  
926  
927  
928  
929  
930  
931  
932  
933  
934  
935  
936  
937  
938  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
970  
971

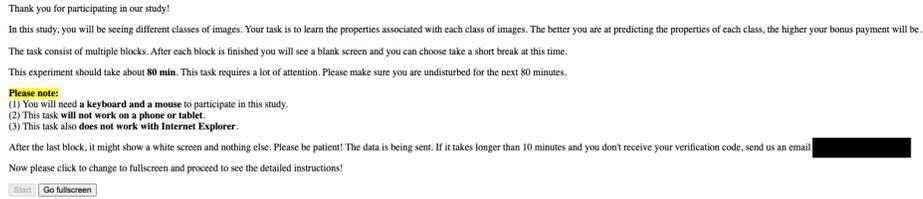


Figure 9: Initial instructions received by participants after the collection of informed consent.

### A.3 DISCRETIZING NETWORK RESPONSES FOR COMPARISON TO HUMANS

**Discretization.** In our task, neural networks produced continuous outputs. This is distinct from human learners who were required to give discrete responses. We now describe the discretization that allows us to compare human and network responses. Fundamentally, we conceptualise inference as a noisy process by which responses are sampled from a distribution over output labels. That is, we treat outputs from our linear network as logits. We first feed network outputs  $\hat{\mathbf{y}}_i$  through a softmax function with temperature 0.2 and subsequently sample three responses without replacement. The procedure maps continuous outputs  $\hat{\mathbf{y}}_i$  to binary responses vectors in  $\{0, 1\}^{N_{out}}$ .

**Expected solutions.** Here we describe the derivation of expected solutions used in Fig. 4, dashed lines. The derivation of these "expected responses" under the sampling procedure allows to make the reliance of network responses on the exact solutions in Section 3 clear.

Consider network outputs  $\hat{\mathbf{y}}_i(t) = \mathbf{W}^2(t)\mathbf{W}^1(t)\mathbf{x}_i$ . We transform these outputs through a softmax function  $\sigma_\beta : \mathbb{R}^{N_{out}} \rightarrow (0, 1)^{N_{out}}$ . Let  $S = \{s_1, s_2, s_3\}$  denote the set of three unique response indices sampled from  $\sigma_\beta(\hat{\mathbf{y}}_i(t))$  without replacement, where  $s_n \in \{1, 2, \dots, N_{out}\}$  for  $n = 1, 2, 3$ , and all  $s_n$  are hence distinct. The probability distribution  $\sigma_\beta(\hat{\mathbf{y}}_i(t))$  is dependent on time  $t$ , therefore denote the produced probability of  $S$  as  $P_t(S)$ . For each of these sets  $S$  we can compute an associated true negative rate for each of the  $k \in \{1, 2, 3\}$  levels in the hierarchy. We denote this random variable as  $X_S^k$ . We can then compute expected solutions to inference behaviour as

$$\mathbb{E}_t[X_S^k] = \sum_{S \subseteq \{1, 2, \dots, m\}, |S|=3} P_t(S) X_S^k \quad (6)$$

### A.4 TRUE NEGATIVE RATES

Here we describe the metric used in the bottom panel of Fig. 3 and in Fig. 4. The metrics effectively describes true negative rates (correct-rejection scores). We use the metric on continuous network responses in  $\mathbb{R}^{N_{out}}$  in Fig. 3. We also use the metric on discretised networks responses in  $\{0, 1\}^{N_{out}}$  and for human responses in  $\{0, 1\}^{N_{out}}$  in Fig. 4.

Given responses  $\hat{\mathbf{y}}$  and target vectors  $\mathbf{y} \in \mathbb{R}^{N_{out}}$  the metric computes the alignment between target and response vectors while only focusing on zero entries in  $\mathbf{y}$ . Furthermore we compute the metric separately for the  $k \in \{1, 2, 3\}$  separate levels of the hierarchy where the entries  $s_k$  and  $e_k$  denote relevant start and end indices of level  $k$  in the vectors  $\hat{\mathbf{y}}$  and  $\mathbf{y}$ . The metric is then computed as

$$f_k^{tn}(\hat{\mathbf{y}}, \mathbf{y}) = \frac{(\mathbf{1}_{N_{out}} - \hat{\mathbf{y}})_{s_k:e_k}^T (\mathbf{1}_{N_{out}} - \mathbf{y})_{s_k:e_k}}{(\mathbf{1}_{N_{out}} - \mathbf{y})_{s_k:e_k} (\mathbf{1}_{N_{out}} - \hat{\mathbf{y}})_{s_k:e_k}}, \quad (7)$$

where  $s_k : e_k$  is a "slicing" notation that takes the subvector between indices  $s_k$  and  $e_k$ .

If for all desired entries of 0 in  $\mathbf{y}$  the vector  $\hat{\mathbf{y}}$  is equal to 0 the metric will be at 1. Correspondingly if entries in  $\hat{\mathbf{y}}$  are larger than zero the metric  $f_k^{tn}(\hat{\mathbf{y}}, \mathbf{y})$  will decrease. Thus, the metric measures wrong beliefs about the presence of target labels across the different levels of the hierarchy.

972 A.5 ADDITIONAL THEORETICAL RESULTS AND PROOFS

973  
974 A.5.1 EQUIVALENCE OF BIAS TERMS

975  
976 In this section, we give more detail on the method used in in Section 3 of how to reformulate a bias  
977 term in terms of the network weights and a constant feature in the input.

978 Consider a network with an explicit input bias term  $\mathbf{b}^1$ ,

979  
980 
$$\hat{\mathbf{y}} = \tilde{\mathbf{W}}^1 \tilde{\mathbf{x}} + \tilde{\mathbf{b}}^1$$

981  
982 This is equivalent to introducing a constant component to the vector  $\mathbf{x}$ ,

983  
984 
$$\tilde{\mathbf{x}} \rightarrow \mathbf{x} := \begin{bmatrix} 1 \\ \tilde{\mathbf{x}} \end{bmatrix},$$

985  
986 and using the network

987  
988 
$$\hat{\mathbf{y}} = \mathbf{W}^1 \mathbf{x},$$

989  
990 as we can write

991  
992 
$$\begin{aligned} 993 (\mathbf{W}^1 \mathbf{x})_m &= \sum_{j=0}^{N_{in}} W_{mj}^1 x_j \\ 994 &= W_{m0}^1 1 + \sum_{j=1}^{N_{in}} W_{mj}^1 x_j \\ 995 &= W_{m0}^1 1 + \sum_{j=0}^{N_{in}-1} \tilde{W}_{mj}^1 \tilde{x}_j \\ 996 &\equiv b_m^1 + \sum_{j=0}^{N_{in}-1} \tilde{W}_{mj}^1 \tilde{x}_j. \end{aligned}$$

997  
998 In order to match a given i.i.d. initialization  $b_m^1 \sim \mathcal{N}(0, \sigma_b^2)$  where  $\sigma_b \neq \sigma_w$ , the component that  
999 needs to be added to  $\tilde{\mathbf{x}}$  to get equivalence needs to be  $\sigma_b/\sigma_w$ .

1000  
1001 A.5.2 LEARNING DYNAMICS FOR BIAS TERMS

1002  
1003 We here derive analytical expressions for the learning speeds of input and output bias terms for a  
1004 two-layer deep linear network discussed in the main text,

1005  
1006 
$$\hat{\mathbf{y}} = \mathbf{W}^2 (\mathbf{W}^1 \mathbf{x} + \mathbf{b}^1) + \mathbf{b}^2.$$

1007  
1008 We decompose  $\mathbf{W}^2 = \mathbf{U} \mathbf{A}^{(2)} \mathbf{R}^{(2)}$  and  $\mathbf{W}^1 = \mathbf{R}^{(1)} \mathbf{A}^{(1)} \mathbf{V}$  by means of a singular value decom-  
1009 position (SVD). We here make the assumption of balancedness  $\mathbf{W}^1(0) \mathbf{W}^{1T}(0) = \mathbf{W}^{2T}(0) \mathbf{W}^2(0)$   
1010 (Braun et al., 2022) at the beginning of training, which implies  $\mathbf{R}^{(2)} \mathbf{S}^{(2)2} \mathbf{R}^{(2)T} = \mathbf{R}^{(1)} \mathbf{S}^{(1)2} \mathbf{R}^{(1)T}$ .  
1011 For clarity, we further assume the simplification

1012  
1013 
$$\mathbf{R}^{(2)T} = \mathbf{R}^{(1)} =: \mathbf{R}, \quad \mathbf{A}^{(2)} = \mathbf{A}^{(1)} =: \sqrt{\mathbf{A}}.$$

1014  
1015 We here just state these relations without further comment to complement the respective derivation  
1016 for the weights in (Saxe et al., 2014). This decomposition then allows to rewrite the gradients.

### Input bias term

$$\begin{aligned}
\tau \frac{d}{dt} \mathbf{b}^1 &= \nabla_{\mathbf{b}^1} \mathcal{L} \\
&= (\mathbf{y} - \hat{\mathbf{y}})^T \mathbf{W}^2 \\
&= (\mathbf{y} - (\mathbf{W}^2 (\mathbf{W}^1 \mathbf{x} + \mathbf{b}^1) + \mathbf{b}^2))^T \mathbf{W}^2 \\
\mathbb{E}_{\mathbf{x}} \rightarrow &(\bar{\mathbf{y}} - \mathbf{W}^2 (\mathbf{W}^1 \bar{\mathbf{x}} + \mathbf{b}^1) - \mathbf{b}^2)^T \mathbf{W}^2 \\
&= (\bar{\mathbf{y}} - \mathbf{UAV}\bar{\mathbf{x}} - \mathbf{U}\sqrt{\mathbf{AR}}\mathbf{b}^1 - \mathbf{b}^2)^T \mathbf{U}\sqrt{\mathbf{AR}}^T \\
&= \bar{\mathbf{y}}^T \mathbf{U}\sqrt{\mathbf{AR}}^T - \bar{\mathbf{x}}^T \mathbf{V}^T \mathbf{AR}^T - \mathbf{b}^{1T} \mathbf{RAR} - \mathbf{b}^{2T} \mathbf{U}\sqrt{\mathbf{AR}}^T \\
&= (\mathbf{Y}\mathbf{1}_N)^T \mathbf{U}\sqrt{\mathbf{AR}}^T - (\mathbf{X}\mathbf{1}_N)^T \mathbf{V}^T \mathbf{AR} - \mathbf{b}^{1T} \mathbf{RAR} - \mathbf{b}^{2T} \mathbf{U}\sqrt{\mathbf{AR}}^T.
\end{aligned}$$

Here, we denoted the expectation over the data samples as  $\mathbb{E}_{\mathbf{x}}$ . Projecting from the right with  $\mathbf{R}_\alpha \in \mathbb{R}^{N_{\text{hidden}}}$  gives

$$\tau \frac{d}{dt} (\mathbf{b}^{1T} \mathbf{R}_\alpha) = \bar{\mathbf{y}}^T \mathbf{U}_\alpha \sqrt{a_\alpha} - \bar{\mathbf{x}}^T \mathbf{V}_\alpha^T a_\alpha - \mathbf{b}^{1T} \mathbf{R}_\alpha a_\alpha - \mathbf{b}^{2T} \mathbf{U}_\alpha \sqrt{a_\alpha}. \quad (8)$$

### Output bias term

$$\begin{aligned}
\tau \frac{d}{dt} \mathbf{b}^2 &= (\mathbf{y} - (\mathbf{W}^2 (\mathbf{W}^1 \mathbf{x} + \mathbf{b}^1) + \mathbf{b}^2)) \\
\mathbb{E}_{\mathbf{x}} \rightarrow &\bar{\mathbf{y}} - \mathbf{W}^2 (\mathbf{W}^1 \bar{\mathbf{x}} + \mathbf{b}^1) - \mathbf{b}^2 \\
&= \bar{\mathbf{y}} - \mathbf{UAV}\bar{\mathbf{x}} - \mathbf{U}\sqrt{\mathbf{AR}}^T \mathbf{b}^1 - \mathbf{b}^2 \\
&= \mathbf{Y}\mathbf{1}_N - \mathbf{UAVX}\mathbf{1}_N - \mathbf{UR}^T \mathbf{b}^1 - \mathbf{b}^2.
\end{aligned} \quad (9)$$

Notably, the derivative in Eq. (8) is proportional to the singular vectors of the weights  $a_\alpha$ , so that its growth is attenuated, analogous to the sigmoidal growth in deep linear networks (Saxe et al., 2014). In contrast, the learning signal  $\frac{d}{dt} \mathbf{b}^2$  in Eq. (9) is not affected by the initialization of the weights and is hence  $\mathcal{O}(1)$  already at the beginning of learning, reminiscent of shallow networks.

#### A.5.3 INTEGRATED FORMULATION OF ARCHITECTURAL BIASES.

In the main text, we have analysed how bias terms on the *input* layer affect the singular value spectrum. Our empirical results in Section 4.1 suggest a more general dynamical bias towards the OCS stemming purely from architectural properties. Here, we use the neural tangent kernel  $\text{NTK}(\mathbf{x}_i, \mathbf{x}_{i'}) = \sum_k \frac{d\hat{\mathbf{y}}_i}{d\theta_k} \frac{d\hat{\mathbf{y}}_{i'}^T}{d\theta_k}$  (Jacot et al., 2018) to directly and comprehensively describe the affected time evolution of the network response  $\frac{d}{dt} \hat{\mathbf{y}}_i = \text{NTK}(\mathbf{y}_i - \hat{\mathbf{y}}_i)$  at the cost of a closed-form solution. Because changes in network outputs are proportional to the NTK it can be viewed as an architecture-induced learning rate (Roberts et al.). For a review and derivation of the NTK, see Appendix A.5.8. For completeness, we now consider a network that contains input  $\mathbf{b}^1$  and output  $\mathbf{b}^2$  bias terms.

**Proposition 3** (NTK of linear networks with bias terms). *Consider a two-layer linear network with input and output-layer bias  $\hat{\mathbf{Y}} = \mathbf{W}^2 (\mathbf{W}^1 \mathbf{X} + \mathbf{b}^1) + \mathbf{b}^2$  in the high-dimensional regime. Furthermore, assume weights are initialized i.i.d.  $W_{ij}^\ell \sim \mathcal{N}(0, \sigma_{\mathbf{W}^\ell}^2 / N_{in}^\ell)$  in each layer. Then, the neural tangent kernel of in early training in expectation  $\mathbb{E}_{\mathbf{W}}$  reads*

$$\text{NTK}(\mathbf{X}, \mathbf{X}) = \sigma_{\mathbf{W}^2}^2 \mathbf{I}_{N_{out}} \otimes \left( 2\mathbf{X}^T \mathbf{X} + \underbrace{\mathbf{1}_N \mathbf{1}_N^T}_{\leftrightarrow \mathbf{b}^1} \right) + \mathbf{1}_{N_{out}} \mathbf{1}_{N_{out}}^T \otimes \underbrace{\mathbf{1}_N \mathbf{1}_N^T}_{\leftrightarrow \mathbf{b}^2}. \quad (10)$$

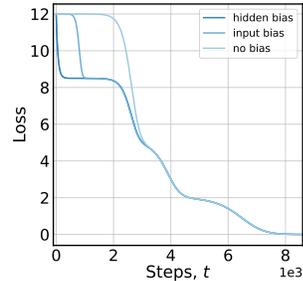


Figure 10: Loss curves for different bias variations.

The tensor product  $\otimes$  separates the components that operate on output and sample space. We briefly review the NTK and derive this expression in the next section. The highlighted terms originate from the bias term  $\frac{d\hat{\mathbf{y}}_i}{d\mathbf{b}} \frac{d\hat{\mathbf{y}}_i^T}{d\mathbf{b}}$  entering the NTK, manifesting in the appearance of the constant mode  $\mathbf{1}_N$ . Importantly, these terms do not scale with the size of the learned bias – they are present even if the bias is initialized at zero. Intuitively, their contribution stems from the architecture’s *potential* to learn a bias, enabling rapid changes in output  $\hat{\mathbf{y}}$ . The NTK also reveals a qualitative difference between input and output bias: Whereas the term that is induced by  $\mathbf{b}^1$  shows attenuated growth due to the multiplication by the weights of initial scale  $\sigma_{\mathbf{W}^2} \ll 1$ , the output bias  $\mathbf{b}^2$  immediately changes the output significantly. Loss curves which demonstrate the effect of different bias terms are displayed in Fig. 10.

## PROOFS

### A.5.4 FEASIBILITY OF CLOSED-FORM SOLUTION

**Proposition 1** (Feasibility of closed-form learning dynamics). *For any input data  $\mathbf{X} \in \mathbb{R}^{N_{in} \times N}$  and output data  $\mathbf{Y} \in \mathbb{R}^{N_{out} \times N}$  it is possible to diagonalize  $\Sigma^x$  by the right singular vectors  $\mathbf{V}$  of  $\Sigma^{yx}$  if  $\mathbf{Y}^T \mathbf{Y}$  and  $\mathbf{X}^T \mathbf{X}$  commute. The converse holds true only if  $\mathbf{X}$  has a left inverse.*

*Proof.* We would like to know when the right singular vectors  $\mathbf{V}$  (denote as  $\mathbf{V}^{yx}$  here for clarity) of  $\Sigma^{yx} = \mathbf{U}^{yx} \mathbf{S}^{yx} \mathbf{V}^{yx}$  match these of  $\Sigma^x = \mathbf{U}^x \mathbf{S}^x \mathbf{V}^x$ . First, to reduce the problem to  $\mathbf{V}^{yx}$ , note that  $\Sigma^{yxT} \Sigma^{yx} = \mathbf{V}^{yx} \mathbf{S}^{yx2} \mathbf{V}^{yx}$ , so that what remains to show is  $[\Sigma^{yxT} \Sigma^{yx}, \Sigma^{xx}] = 0$ , where  $[\mathbf{A}, \mathbf{B}] := \mathbf{A}\mathbf{B} - \mathbf{B}\mathbf{A}$  denotes the commutator between two matrices  $\mathbf{A}$  and  $\mathbf{B}$ . We compute the two terms as

$$\begin{aligned} \Sigma^{yxT} \Sigma^{yx} \Sigma^{xx} &= \mathbf{X}\mathbf{Y}^T \mathbf{Y}\mathbf{X}^T \mathbf{X}\mathbf{X}^T \\ \Sigma^{xxT} \Sigma^{yxT} \Sigma^{yx} &= \mathbf{X}\mathbf{X}^T \mathbf{X}\mathbf{Y}^T \mathbf{Y}\mathbf{X}^T \end{aligned}$$

The commutator vanishes if these terms match, which happens for the simpler equality

$$\mathbf{Y}^T \mathbf{Y}\mathbf{X}^T \mathbf{X} = \mathbf{X}^T \mathbf{X}\mathbf{Y}^T \mathbf{Y},$$

or  $[\mathbf{Y}^T \mathbf{Y}, \mathbf{X}^T \mathbf{X}] = 0$ . The converse follows only if the transformation  $\mathbf{X} \dots \mathbf{X}^T$  in the former equation is invertible, which is the case if a left inverse  $\mathbf{X}^{-1} \mathbf{X} = \mathbf{I}_{N_{in}}$  exists.

### A.5.5 OCS AND SHARED PROPERTIES CORRESPOND TO EACH OTHER

We here link the OCS and shared properties stand in close relation, as the eigenvector  $\mathbf{1}_N$  represents properties that are shared across all data samples.

**Proposition 2** (The OCS is linked to shared properties). *If  $\mathbf{1}_N$  is an eigenvector to the similarity matrix  $\mathbf{X}^T \mathbf{X} \in \mathbb{R}^{N \times N}$ , then the sample-average  $\bar{\mathbf{x}} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i$  will be an eigenvector to the correlation matrix  $\mathbf{X}\mathbf{X}^T \in \mathbb{R}^{N_{in} \times N_{in}}$  with identical eigenvalue  $\lambda$ . An analogous statement applies for  $\mathbf{Y}^T \mathbf{Y}$  and  $\mathbf{Y}\mathbf{Y}^T$ . The converse does not hold true in general.*

*Proof.*

$$\mathbf{X}\mathbf{X}^T \bar{\mathbf{x}} = \left( \mathbf{X}\mathbf{X}^T \right) \frac{1}{N} \mathbf{X} \mathbf{1}_N = \frac{1}{N} \mathbf{X} (\mathbf{X}^T \mathbf{X}) \mathbf{1}_N = \frac{1}{N} \mathbf{X} \lambda \mathbf{1}_N = \lambda \frac{1}{N} \mathbf{X} \mathbf{1}_N = \lambda \bar{\mathbf{x}}.$$

### A.5.6 CONSTANT DATA MODE $\mathbf{1}_N$ IS RELATED TO SYMMETRY IN THE DATA

This section gives proof sketches based on symmetry in the dataset that are sufficient to make  $\mathbf{1}_N$  an eigenvector to  $\mathbf{X}^T \mathbf{X}$  and  $\mathbf{Y}^T \mathbf{Y}$ , and in particular hold for the dataset that we are considering. We anticipate that it is possible to formulate these statements in a more universal way by fully leveraging the cited literature.

The assumptions on symmetry should intuitively at least hold in an approximate manner for many datasets, we expect that they indeed are the reason why we observe a prevalence of  $\mathbf{1}_N$ , although they are not a necessary condition.

**Continuously supported data**  $\mathbf{x} \in \mathbb{R}^{N_{in}}$

**Proposition 4** (Continuous symmetry induces  $\mathbf{1}_N$ ). *If the pairwise correlations  $\mathbf{y}_i^T \mathbf{y}_{i'}$  in a dataset are rotationally symmetric, its similarity matrix  $\mathbf{Y}^T \mathbf{Y}$  has eigenvector  $\mathbf{1}_N$ . Note that this is a weaker assumption than the data itself being symmetric.*

*Proof.* We assume that  $\mathbf{Y}, \mathbf{X}$  have been sampled from a ground truth data distribution  $p(\mathbf{y}, \mathbf{x})$ . If  $p(\mathbf{y})$  is rotationally symmetric and  $\mathbf{X}$  is comprised of samples  $\mathbf{x}$  that are uniformly distributed on the hypersphere, we can introduce the kernel function  $\mathbf{y}_{\mathbf{x}_i}^T \mathbf{y}_{\mathbf{x}_{i'}} = k(\mathbf{x}_i, \mathbf{x}_{i'}) = k(\mathbf{R}_l \mathbf{x}_i, \mathbf{R}_l \mathbf{x}_{i'}) = k(\mathbf{x}_i^T \mathbf{x}_{i'})$  for any  $R_l$  that is a representation of the group of rotations  $G = \text{SO}(N_{in})$  that faithfully acts on the “subsamped” hypersphere  $\mathbf{X}$  comprised of vectors  $\mathbf{x} \in \mathbb{R}^{N_{in}}$ . It therefore only depends on the pairwise input similarity (hence sometimes called dot-product kernel). It follows that for all vectors  $\mathbf{v}(\mathbf{X}) \in \mathbb{R}^N$  that are evaluations of the functions of the sample points  $\mathbf{X}$

$$\mathbf{Y}^T \mathbf{Y} \mathbf{v} = k(\mathbf{X}^T \mathbf{X}) \mathbf{v} = k((\mathbf{R}_l \mathbf{X})^T (\mathbf{R}_l \mathbf{X})) \mathbf{v} = \mathbf{R}_k^T k(\mathbf{X}^T \mathbf{X}) \mathbf{R}_l \mathbf{v} \Leftrightarrow [\mathbf{Y}^T \mathbf{Y}, \mathbf{R}_l] = 0,$$

where  $[\mathbf{A}, \mathbf{B}] =: \mathbf{AB} - \mathbf{BA}$  is the commutator between two matrices.

It follows that we must have for all rotations  $\mathbf{R}_l$

$$\mathbf{R}_l (\mathbf{Y}^T \mathbf{Y} \mathbf{1}_N) = \mathbf{Y}^T \mathbf{Y} \mathbf{R}_l \mathbf{1}_N = \mathbf{Y}^T \mathbf{Y} \lambda_{\mathbf{R}_l} \mathbf{1}_N = \lambda_{\mathbf{R}_l} \mathbf{Y}^T \mathbf{Y} \mathbf{1}_N$$

with eigenvalue  $\lambda_{\mathbf{R}_l} = 1$ .

meaning that  $\mathbf{Y}^T \mathbf{Y} \mathbf{1}_N$  is an eigenvector to *all*  $\mathbf{R}_l$ . This can only be the case if  $\mathbf{Y}^T \mathbf{Y} \mathbf{1}_N \propto \mathbf{1}_N$ , as this is the only vector of values on the sphere that is invariant under any rotations.

We point out that it can be shown more generally with tools from functional analysis that the full spectrum of this kernel operator  $k$  are the spherical harmonics if the data measure  $p(\mathbf{x})$  is spherically symmetric (Hecke, 1917), see (Dutordoir et al., 2020) for a modern presentation with tools from calculus. As the first harmonic  $\mathcal{Y}_{l=0, m=0}(\mathbf{x})$  is constant, it follows that also the constant function  $\mathbf{1}(\mathbf{x}) \equiv 1$  is an eigenfunction when drawing a finite set of samples from this kernel.

**Data on a graph**  $\mathbf{x} \in \mathbb{R}^{N_{in}}$  We here prove that the former statement holds for the hierarchical dataset that is discussed in the main text, i.e. that  $\mathbf{1}_N$  is an eigenvector to  $\mathbf{Y}^T \mathbf{Y}$ .

First, note that it is easy to convince oneself of this by writing down the matrices explicitly: Then, as the rows are just permutations of one another,  $\mathbf{1}_N$  is immediately identified as an eigenvector, because  $\sum_{i'} \mathbf{Y}_i^T \mathbf{Y}_{i'} \mathbf{1}_{i'} = \mathbf{Y}_i^T (\sum_{i'} \mathbf{Y}_{i'})$  will then not depend on  $i$  and hence be proportional to  $\mathbf{1}_N$ .

To connect with the former symmetry-based argument Appendix A.5.6, we here however give a proof that is based on the symmetry in the data:

**Proposition 5** (Discrete symmetry induces  $\mathbf{1}_N$ ). *Consider a connected Cayley tree graph with adjacency matrix  $\mathbf{A}$  and nodes  $\mathbf{x}_i$ . Furthermore, let  $\mathbf{R}_l \in G$  be an element of a faithful representation of the symmetry group  $G$  that acts on the graph nodes  $\mathbf{v}$ , i.e. that leaves its adjacency matrix invariant,  $[\mathbf{R}_l, \mathbf{A}] = 0 \forall \mathbf{R}_l$ .*

*If  $\mathbf{Y}$  are labels associated with the leaf nodes  $\mathbf{X}$  (the outermost generation of the graph, see (Erzan and Tuncer, 2020)) and there exists a similarity function  $k$  such that  $\mathbf{y}_{\mathbf{x}_i}^T \mathbf{y}_{\mathbf{x}_{i'}} = \mathbf{y}_{\mathbf{R}_l \mathbf{x}_i}^T \mathbf{y}_{\mathbf{R}_l \mathbf{x}_{i'}} = k(d(\mathbf{x}_i, \mathbf{x}_{i'})) \forall \mathbf{R}_l$  where  $d$  is the geodesic distance on the graph,  $\mathbf{1}_N$  will be an eigenvector of  $\mathbf{Y}^T \mathbf{Y}$ .*

1188 *Proof.* From the symmetry assumption on the labels, we again have for any vector  $\mathbf{v}$  of node loadings  
 1189  $[\mathbf{R}_l, \mathbf{Y}^T \mathbf{Y}] \mathbf{v} = 0 \forall \mathbf{R}_l \in G$ . From this, we find that

$$1190 \mathbf{R}_l (\mathbf{Y}^T \mathbf{Y} \mathbf{1}_N) = \mathbf{Y}^T \mathbf{Y} \mathbf{R}_l \mathbf{1}_N = \mathbf{Y}^T \mathbf{Y} \lambda_{\mathbf{R}_l} \mathbf{1}_N = \lambda_{\mathbf{R}_l} \mathbf{Y}^T \mathbf{Y} \mathbf{1}_N \forall \mathbf{R}_l$$

1192 This shows that  $\mathbf{Y}^T \mathbf{Y} \mathbf{1}_N$  is an eigenvector of  $\mathbf{R}_l$  with eigenvalue  $\lambda_{\mathbf{R}_l} = 1$  for any element of the  
 1193 symmetry group. The only vector  $\mathbf{v}$  that is invariant under *all* symmetry operations of the graph is  
 1194 the constant vector  $\mathbf{1}_N$ .

1196 We briefly point out the rich literature on spectral graph theory (for example (Brouwer and Haemers,  
 1197 2011; Erzan and Tuncer, 2020)) that might allow making statements about the nature of the eigen-  
 1198 values and other eigenvectors as a function of the graph topology. We expect that this is possible  
 1199 because the literature in the continuous case discussed in the next paragraph bases their arguments  
 1200 on the Laplacian on the sphere, an operator that can be extended to graphs as well. We leave these  
 1201 exploration for future work.

1202 **Corollary 2.** *Because  $k(\mathbf{X}^T \mathbf{X}) := \mathbf{X}^T \mathbf{X}$  defines a particular case of input-output similarity map-  
 1203 ping,  $\mathbf{1}_N$  is also an eigenvector to  $\mathbf{X}^T \mathbf{X}$  under the former assumptions of uniform data distribution.*

#### 1204 A.5.7 CONSTANT DATA MODE $\mathbf{1}_N$ IS THE LEADING EIGENVECTOR

1205 Here, we prove that the constant eigenvector  $\mathbf{1}_N$  which is responsible for the OCS solution is  
 1206 associated with the *leading* eigenvalue of the input-output correlation matrix and hence drives early  
 1207 learning.

1208 **Theorem 1** (Early learning is biased by the OCS mode). *If  $\mathbf{1}_N$  is a joint non-degenerate eigenvector  
 1209 to positive input and output similarity matrices  $\mathbf{X}^T \mathbf{X}$  and  $\mathbf{Y}^T \mathbf{Y}$ , the OCS mode  $s_{ocs} \bar{\mathbf{y}} \bar{\mathbf{x}}^T$  will have  
 1210 leading spectral weight  $s_0 \equiv s_{ocs}$  in the SVD of the input-output correlation matrix  $\Sigma^{yx}$ .*

1211 *Proof.* Let  $\mathbf{1}_N$  be an eigenvector to both similarity matrices  $\mathbf{X}^T \mathbf{X}$  and  $\mathbf{Y}^T \mathbf{Y}$  associated with  
 1212 eigenvalue  $\tilde{\lambda}$ . Moreover, let  $\mathbf{X}^T \mathbf{X}$  and  $\mathbf{Y}^T \mathbf{Y}$  have positive entries. Then, the Perron-Frobenius  
 1213 theorem (Perron, 1907) guarantees that  $\tilde{\lambda}$  is indeed the leading eigenvalue to  $\mathbf{Y}^T \mathbf{Y}$ ,  $\tilde{\lambda} \equiv \lambda_0 = s_0^2$ .

1214 By Proposition 2,  $\bar{\mathbf{x}}$  and  $\bar{\mathbf{y}}$  are now also the leading eigenvectors for  $\mathbf{X} \mathbf{X}^T$  and  $\mathbf{Y} \mathbf{Y}^T$ . Because the  
 1215 eigenvectors of  $\mathbf{Y} \mathbf{Y}^T$  and  $\mathbf{X} \mathbf{X}^T$  are the left and right singular vectors of  $\Sigma^{yx}$ , respectively, with the  
 1216 eigenvalues being the squares of the singular values, it follows that

$$1217 s_0 \mathbf{u}_0 \mathbf{v}_0^T = \sqrt{\lambda_0} \bar{\mathbf{y}} \bar{\mathbf{x}}^T.$$

#### 1223 A.5.8 NEURAL TANGENT KERNEL

1224 In this section, we review the neural tangent kernel (NTK). This object is useful because it directly  
 1225 describes the learning dynamics in output space  $\hat{y}$  (Jacot et al., 2018; Roberts et al.) as we briefly  
 1226 demonstrate here. We then calculate the NTK for our specific architecture to yield ?? in the main  
 1227 text. The following makes use of Einstein summation convention.

1228 For a vector-valued model  $\hat{\mathbf{y}}(\mathbf{x}) \in \mathbb{R}^{N_{out}}$  parametrized by a parameter vector  $\theta$ , the evaluation on  
 1229 sample  $\mathbf{x}_i$  from training data at  $\mathbf{x}_{i'}$  evolves as

$$1230 \tau \frac{d}{dt} y_m(\mathbf{x}_i) = \sum_k \frac{dy_m(\mathbf{x}_i)}{d\theta^k} \frac{d\theta^k}{dt} \quad (11)$$

$$1231 = -\eta \sum_k \frac{dy_m(\mathbf{x}_i)}{d\theta^k} \frac{d\mathcal{L}}{d\theta^k} \quad (12)$$

$$1232 = -\eta \left[ \sum_k \frac{dy_m(\mathbf{x}_i)}{d\theta^k} \frac{dy_{m'}(\mathbf{x}_{i'})}{d\theta^k} \right] \frac{d\mathcal{L}}{dy_{m'}}(\mathbf{x}_{i'}) \quad (13)$$

$$1233 =: -\eta \text{NTK}_{mm'}(\mathbf{x}_i, \mathbf{x}_{i'}) (y_{m'}(\mathbf{x}_{i'}) - \hat{y}_{m'}(\mathbf{x}_{i'})), \quad (14)$$

1242 where we used the chain rule and that the parameters update according to gradient descent with  
1243 learning rate  $\eta$ ,  $\frac{d\theta^k}{dt} = -\eta \frac{d\mathcal{L}}{d\theta^k}$ . The last line has defined the NTK. We set  $\eta = 1$  in the main text for  
1244 simplicity, as it does not change trajectory and thereby convergence in the case of gradient flow. In  
1245 addition, we evaluated  $\frac{d\mathcal{L}}{dy_{m'}}(\mathbf{x}_{i'})$  for the case of MSE loss  $\mathcal{L}(\mathbf{x}_{i'}) = 1/2 \sum_{m'} (y_{m'}(\mathbf{x}_{i'}) - \hat{y}_{m'}(\mathbf{x}_{i'}))^2$ .  
1246 The last line of Eq. (11) reveals that the NTK acts as an effective learning rate, as noted by Roberts  
1247 et al..  
1248

1249 We here consider a two-layer linear architecture  $\hat{Y}_m^i(\mathbf{X}) = W_{mk}^2 \left( W_{kj}^1 X_j^i + b_k^1 \right) + b_m^2$  where  
1250 we adopt Einstein summation convention over repeated indices. The parameters are  $\theta^k \in$   
1251  $\{\mathbf{W}^2, \mathbf{W}^1, \mathbf{b}^1, \mathbf{b}^2\}$ . Herein,  $m$  indexes output features and  $i$  indexes data samples. The non-zero  
1252 gradients are  
1253

$$\begin{aligned} 1254 \quad \frac{d\hat{Y}_m^i}{dW_{mk}^2} &= W_{kj}^1 X_j^i + b_k^1 \\ 1255 \quad \frac{d\hat{Y}_m^i}{db_m^2} &= 1_m \\ 1256 \quad \frac{d\hat{Y}_m^i}{dW_{kj}^1} &= W_{mk}^2 X_j^i \\ 1257 \quad \frac{d\hat{Y}_m^i}{db_k^1} &= W_{mk}^2 \mathbf{1}_k. \end{aligned}$$

1264 Inserting this into Eq. (11), we get  
1265

$$\begin{aligned} 1266 \quad \text{NTK}_{m_1 m_2}(X_j^{i_1}, X_j^{i_2}) &= I_{m_1 m_2} \left( X_{j'}^{i_1} W_{j'k}^1 W_{kj''}^1 X_{j''}^{i_2} + b_k^1 b_k^1 \right) \\ 1267 \quad &+ \mathbf{1}_{m_1} \mathbf{1}_{m_2} \\ 1268 \quad &+ W_{m_1 k}^2 W_{km_2}^{2T} X_j^{i_1} X_j^{i_2} \\ 1269 \quad &+ W_{m_1 k}^2 \mathbf{1}_k \mathbf{1}_k W_{km_2}^2. \end{aligned}$$

1272 or in matrix notation, collecting similar terms  
1273

$$\begin{aligned} 1274 \quad \text{NTK}(\mathbf{X}, \mathbf{X}) &= \mathbf{I}_{N_{out}} \otimes \mathbf{X}^T \mathbf{W}^{1T} \mathbf{W}^1 \mathbf{X} + \mathbf{b}^{1T} \mathbf{b}^1 \\ 1275 \quad &+ \mathbf{1}\mathbf{1}^T \otimes \underbrace{\mathbf{1}\mathbf{1}^T}_{\leftrightarrow \mathbf{b}^2} \\ 1276 \quad &+ \mathbf{W}^2 \mathbf{W}^{2T} \otimes \left( \mathbf{X}^T \mathbf{X} + \underbrace{\mathbf{1}\mathbf{1}^T}_{\leftrightarrow \mathbf{b}^1} \right) \\ 1277 \quad &\in \mathbb{R}^{N_{out} \times N_{out}} \otimes \mathbb{R}^{N \times N}, \end{aligned}$$

1282 where the left hand side operator in the tensor product  $\otimes$  is acting in output space  $m_1 m_2$ , whereas  
1283 the right hand side operator acts in pattern space  $i_1 i_2$ . The notation  $\leftrightarrow \mathbf{b}$  indicates that a term is due  
1284 to the bias term. To illustrate this, the NTK acts on the set of labels  $\mathbf{Y} \in \mathbb{R}^{N_{in} \times N}$  as follows:  
1285

$$1286 \quad (\text{NTK}(\mathbf{X}, \mathbf{X}) \mathbf{Y})_i^m = \sum_{m'}^{N_{out}} \sum_{i'}^N \text{NTK}(X_i, X_{i'})^{mm'} Y_{i'}^{m'}. \quad (15)$$

1288 For simplicity, we approximate  $\mathbf{W}^2(0) \mathbf{W}^{2T}(0) = \sigma_{\mathbf{W}}^2 I_{N_{out}}$  and  $\mathbf{W}^{1T}(0) \mathbf{W}^1(0) = \sigma_{\mathbf{W}}^2 I_{N_{in}}$ , which  
1289 approximately holds for initialization  
1290

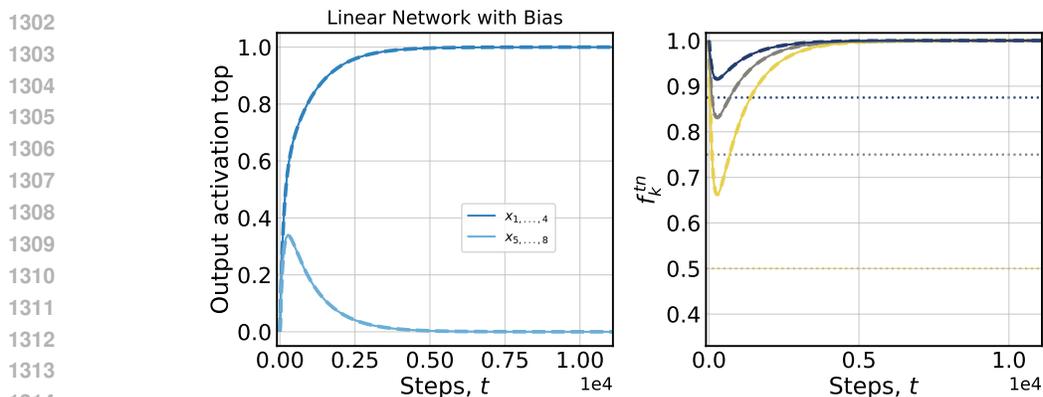
$$1291 \quad \mathbf{W}^1(0) \sim \mathcal{N}(0, \sigma_{\mathbf{W}}^2 / N_{hid}), \quad \mathbf{W}^2(0) \sim \mathcal{N}(0, \sigma_{\mathbf{W}}^2 / N_{hid}), \quad \mathbf{b}^1 = 0, \quad \mathbf{b}^2 = 0$$

1292 where  $N_{hid}$  is the size of the hidden layer and both  $N_{in}$  and  $N_{hid}$  are large. This leaves  
1293

$$1294 \quad \text{NTK}(\mathbf{X}, \mathbf{X}) = \sigma_{\mathbf{W}}^2 \mathbf{I}_{N_{out}} \otimes \left( 2\mathbf{X}^T \mathbf{X} + \underbrace{\mathbf{1}\mathbf{1}^T}_{\leftrightarrow \mathbf{b}^1} \right) + \mathbf{1}\mathbf{1}^T \otimes \underbrace{\mathbf{1}\mathbf{1}^T}_{\leftrightarrow \mathbf{b}^2}.$$

1296 A.6 SHALLOW NETWORK OCS LEARNING

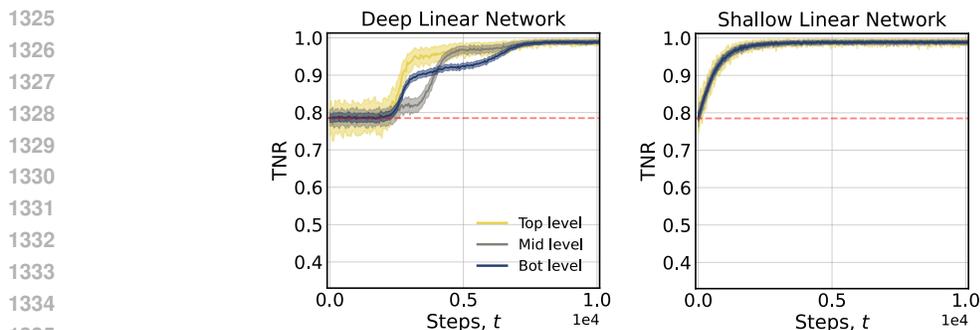
1297  
 1298 In this brief section we show the OCS signatures of shallow networks with bias terms. The result is  
 1299 displayed in Fig. 11. We see similar behavioural signatures to deep linear networks. However, the  
 1300 tendency to the OCS is more transient.  
 1301



1315 Figure 11: Early learning in shallow networks with bias terms approaches the OCS.

1316  
 1317  
 1318 A.7 TRUE NEGATIVE RATES IN LINEAR NETWORKS WITHOUT BIAS TERMS

1319  
 1320 In this short section we provide a supplemental figure relevant for our results in Section 5: We train  
 1321 deep and shallow linear networks *without* bias terms. The learning setting and computation of metrics  
 1322 are equivalent to results in Fig. 4. We display the result in Fig. 12. While networks learn the task,  
 1323 early, response biases are fully absent in these models.  
 1324



1337 Figure 12: True negative rates for linear networks *without* bias terms. We do not see characteristic  
 1338 response patterns observed in Fig. 4.  
 1339

1340  
 1341 A.8 CNN DATASETS AND HYPERPARAMETERS

1342 **Datasets used.** We used and adapted different image datasets for our experiments with CNNs. While  
 1343 the main text focused on results obtained with a variant of MNIST we report further experiments we  
 1344 conducted to highlight the generality of early OCS learning.  
 1345

- 1346 1. **Hierarchical MNIST.** We used the default ten digit classes provided by MNIST. We  
 1347 then sampled 8 classes randomly and replaced the default one-hot labels corresponding to  
 1348 each class  $i$  with the hierarchical, “three-hot” labels  $y_i$  as seen in Fig. 1. E.g., all images  
 1349 corresponding to MNIST digit “1” might be assigned some random “three-hot” output vector  
 $y_i$ .

- 
- 1350 2. **Hierarchical CIFAR-10.** We applied the same procedure and randomly sampled eight  
1351 classes from CIFAR-10 and replaced one-hot labels as for the hierarchical MNIST.
  - 1352 3. **Imbalanced-binary-MNIST.** While not described in the main text we also report results in  
1353 a setting with standard one-hot target vectors. We randomly sample two MNIST classes for  
1354 training. To assess the impact of the OCS in early learning we introduced class imbalanced  
1355 by oversampling one of the two classes by a factor of two.
  - 1356 4. **Standard CelebA.** We perform experiments on CelebA’s face attribute detection task. The  
1357 task offers a natural testbed for early learning of the OCS as face attribute target labels form  
1358 a non-uniform distribution as seen in Fig. 15, bottom. We also normalised images in the  
1359 dataset before training.

1361 **Model details.** We trained a custom CNN with 3 convolutional layers (layer 1: 32 filters of size  $5 \times 5$ ;  
1362 layer 2: 64 filters of size  $3 \times 3$ ; layer 3: 96 filters of size  $3 \times 3$ ), followed by 2 fully connected layers of  
1363 sizes 512 and 256. Activation functions for all layers were chosen as ReLUs. The final layer of the  
1364 model did not contain an activation function when training with squared error loss. In experiments  
1365 with the class imbalanced-binary-MNIST and cross-entropy loss the final layer contained a softmax  
1366 function as non-linearity. For experiments on CelebA the final layer contained sigmoid activation  
1367 functions and we trained with a binary cross-entropy loss over all 40 labels.

1368 **Training details.** For our results on hierarchical MNIST we train models with minibatch SGD with  
1369 a batch size of 16 and with a relatively small step size of  $1e-4$  to examine the early learning phase.  
1370 For all experiments we used Xavier uniform initialisation (Glorot and Bengio, 2010). Whenever  
1371 we use bias terms in the model we initialize these as 0 in line with common practice. For our main  
1372 experiments we train models using a simple squared error loss function. However, to demonstrate  
1373 generality we repeat experiments for the case of class imbalance using a cross-entropy loss and  
1374 binary cross-entropy in the case of CelebA. All experiments are repeated 10 times with different  
1375 random seeds with the exception of CelebA where we used 5 different random seeds, we provide  
1376 standard errors in all figures (shaded regions). For experiments on the hierarchical CIFAR-10, the  
1377 class imbalanced MNIST, and CelebA we kept all parameters as above but we increase step size to  
1378  $1e-3$ . We trained CNN models on an internal cluster on a single RTX 5000 GPU. Runs took less than  
1379 one hour to complete.

## 1380 A.9 ADDITIONAL EXPERIMENTAL RESULTS

1381  
1382 To understand the generality of OCS learning we plot the results of experiments examining early  
1383 learning of the OCS in these models. We mostly restrict ourselves to plots as seen in Fig. 5 as we  
1384 deem these figures most instructive.

1385 **Hierarchical CIFAR-10.** We train on a hierarchical version of CIFAR-10. Where we randomly  
1386 sample 8 classes from MNIST and replace target labels by hierarchical vectors as in Section 4. We  
1387 find the key signatures of early OCS learning: We find early indifference, the reversion of performance  
1388 metrics to the OCS, and a small initial distance of average response to the OCS solution. The  
1389 results mirror behaviour on the hierarchical MNIST shown in Fig. 3 and Fig. 5.

1390 **Class-imbalance MNIST.** We train on an imbalanced MNIST task as described in Appendix A.8.  
1391 We plot the results for training with squared error and cross-entropy loss functions in Fig. 14. Both  
1392 settings show reversion to the OCS. Note that average model outputs in the case of the cross-entropy  
1393 loss start relatively close to outputs expected under the OCS. Despite this proximity the model is still  
1394 driven towards the OCS solution. The results on this imbalanced case highlight potential fairness  
1395 implications. Given that network have been found to revert to the OCS when generalising (Kang  
1396 et al., 2024), early learning in the OCS setting can transiently, but significantly disadvantage minority  
1397 classes. We further highlight this point in a second solvable case of linear networks with bias terms  
1398 in Appendix A.10.

1399 **Standard CelebA.** We show distance from the OCS for the CelebA face attribute detection task  
1400 in Fig. 15, top. CelebA provides a useful test for our hypothesis as attribute labels display natural  
1401 imbalances. We highlight the strong non-uniformity of the majority attribute labels in Fig. 15,  
1402 bottom. We again train networks in two variants: one with squared-error loss and one with binary  
1403 cross-entropy loss applied over all 40 face attributes. With both loss functions network responses are  
driven towards the OCS in early learning. This case further highlights the generality of early OCS

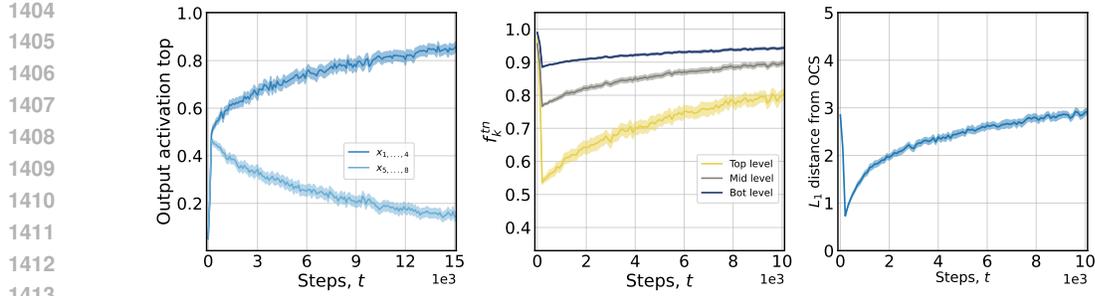


Figure 13: Early OCS learning CNNs trained on hierarchical CIFAR-10. *left*: Network outputs for a single output unit in response to all inputs  $x_i$ . *Centre*: Performance metrics  $f^{tn}$  (Appendix A.4). *Right*: Mean distance of network responses from OCS. Averages taken over every 10 batches for plotting.

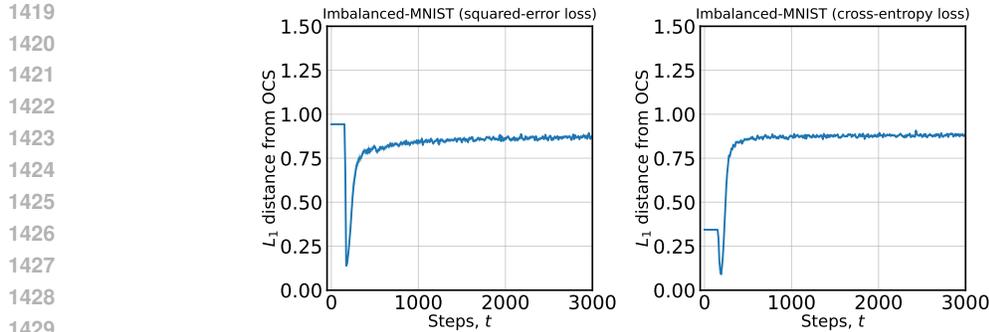


Figure 14: Mean distance of network responses from OCS in CNNs trained on the imbalanced MNIST task. Averages taken over each batch.

learning. OCS learning might be especially undesirable in this setting for fairness reasons as the model will be overly liberal in the prediction of the most common face attributes.

#### A.10 LINEAR NETWORKS UNDER CLASS IMBALANCE

In this section, we describe a second case of a solvable linear network with bias terms. Our dataset consists of two examples where one example appears twice as frequently. We show the data used on the right side of Fig. 16. The minority class has two identifying labels, while this construction appears artificial, it allows for the application of Proposition 1 and solutions to learning dynamics from Section 2 apply.

The case is of particular practical relevance as it illustrates the impact of early OCS learning under class imbalance, a common problem in machine learning where datasets are often naturally imbalanced (Feldman, 2020; Van Horn and Perona, 2017). In practice, these settings are often addressed through oversampling of minority classes (Haibo He and Garcia, 2009; Huang et al., 2016). Empirical work by Ye et al. (2021) documented that neural networks initially fail to learn information about the minority class while classifying most minority examples as belonging to the majority class. Subsequent theoretical work by Francazi et al. (2023) demonstrated that the phenomenon is caused by competition between the optimisation of different classes.

Our work adds to this literature by providing dynamics in a case of gradient-based learning under class imbalance learning that is exactly solvable. Our exact solutions highlight the potential role of early OCS learning in the initial failure to learn about minority classes. The OCS solution substantially biases early predictions towards the majority class as seen in Fig. 16, centre. The results also can be understood as solvable analogous to early reversion to the OCS seen in the Imbalanced-binary-MNIST setting in Fig. 14. The results highlight the potential fairness implications of early OCS learning as the learning phase systematically biases the model against the minority classes.

1458  
 1459  
 1460  
 1461  
 1462  
 1463  
 1464  
 1465  
 1466  
 1467  
 1468  
 1469  
 1470  
 1471  
 1472  
 1473  
 1474  
 1475  
 1476  
 1477  
 1478  
 1479  
 1480  
 1481  
 1482  
 1483  
 1484  
 1485  
 1486  
 1487  
 1488  
 1489  
 1490  
 1491  
 1492  
 1493  
 1494  
 1495  
 1496  
 1497  
 1498  
 1499  
 1500  
 1501  
 1502  
 1503  
 1504  
 1505  
 1506  
 1507  
 1508  
 1509  
 1510  
 1511

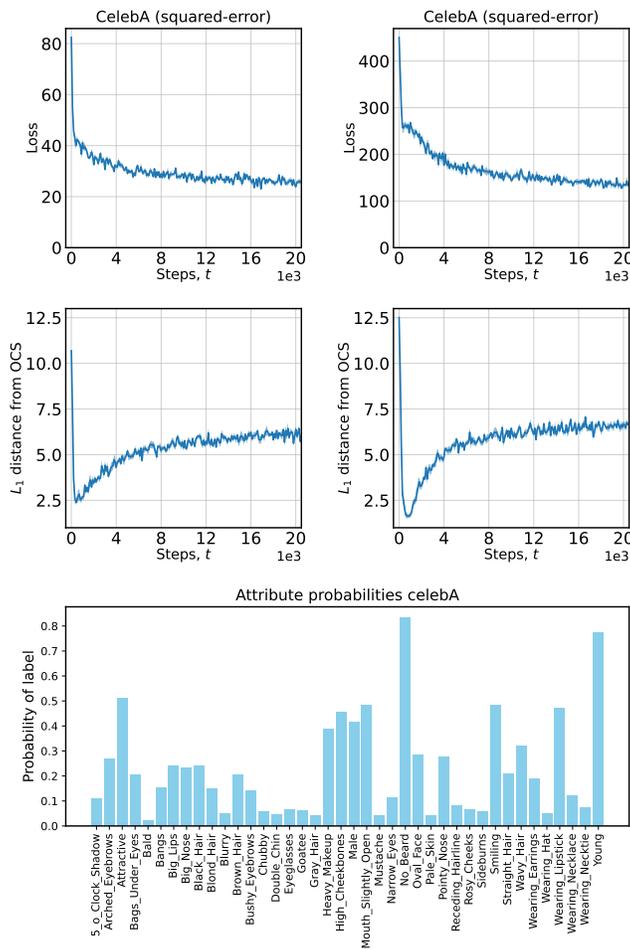


Figure 15: *Top*: Mean distance of network responses from the OCS in CNNs trained on the CelebA face attribute prediction task. *Bottom*: Marginal probabilities of CelebA face attributes.

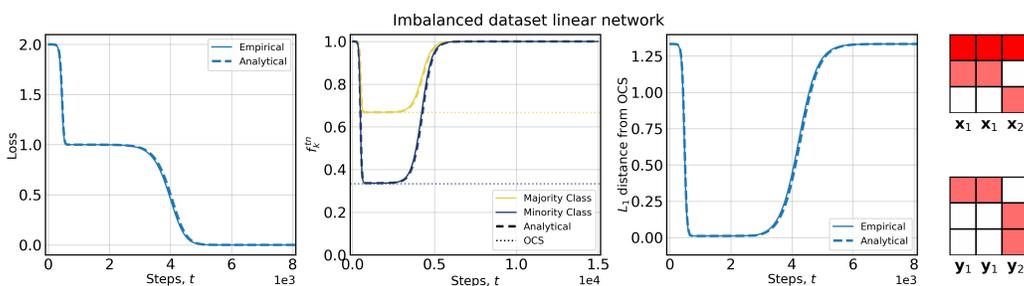


Figure 16: Early learning of the OCS in linear networks under class imbalance.

A.11 OCS-LEARNING IN LINEAR NETWORKS WITH INPUT CORRELATIONS.

In this section, we demonstrate how input correlations can drive OCS learning in the absence of bias terms in linear networks. Specifically we highlight how OCS learning can emerge if  $\mathbf{1}_N$  is an eigenvector of the data similarity matrix  $\mathbf{X}^T \mathbf{X}$ . Note that the network contains no input correlations. In the bottom row of the Fig. 17, we can see that the first SVD mode  $\mathbf{u}_1 \mathbf{v}_1^T$  is indeed exactly equivalent to the OCS mode, i.e.  $\bar{\mathbf{y}} \bar{\mathbf{x}}^T$ . The right panel highlights how the network is driven towards the OCS up until the time-point when the second effective singular value  $a_2(t)$  (which is quite close in time) is learned.

1512  
 1513  
 1514  
 1515  
 1516  
 1517  
 1518  
 1519  
 1520  
 1521  
 1522  
 1523  
 1524  
 1525  
 1526  
 1527  
 1528  
 1529  
 1530  
 1531  
 1532  
 1533  
 1534  
 1535  
 1536  
 1537  
 1538  
 1539  
 1540  
 1541  
 1542  
 1543  
 1544  
 1545  
 1546  
 1547  
 1548  
 1549  
 1550  
 1551  
 1552  
 1553  
 1554  
 1555  
 1556  
 1557  
 1558  
 1559  
 1560  
 1561  
 1562  
 1563  
 1564  
 1565

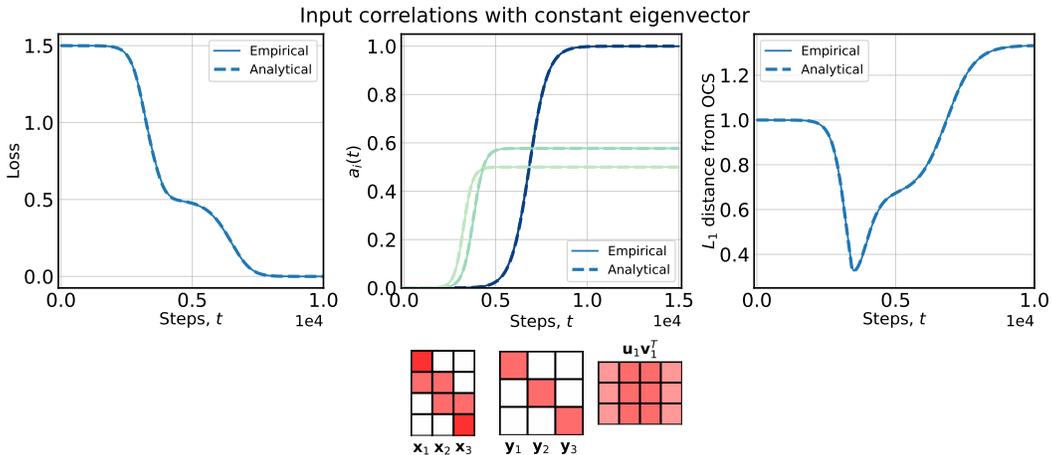


Figure 17: Early learning of the OCS in linear networks with input correlations but in the absence of bias terms.

### A.12 RELATION TO IMBALANCED MULTI-LABEL LEARNING

Given the hierarchical structure of labels used in the majority of our experiments we also see some general connections of our work to problems in the domain of imbalanced multi-label learning. Multi-label learning deals with learning problems in which a single input example is associated with multiple output labels simultaneously. In these settings class imbalance is a key challenge that frequently hinders good performance of models (Liu et al., 2020; Charte et al., 2015; Pham et al., 2021; Liu et al., 2022). Similar to standard classification problems model biases are frequently addressed through adjustments to the models loss function via selective reweighing (Cui et al., 2019) or through sampling based methods which selectively over- or under-sample particular labels (Charte et al., 2015) or via both methods (Pham et al., 2021). Our results on the hierarchical learning task and on the problem of class imbalanced learning in Appendix A.10 might hint at OCS learning as a potential contributor to problems observed in multi-label learning as the imbalanced distribution of output labels might drive learning to undesirable solutions.