



Rebuttal Figure 1: (a) Learning trajectories for two algorithms: $\gamma = 1$ vs. $\gamma = -1$. (b) Generalizing task to pixel inputs. Top row: texture stimuli (A to H). Bottom row: Covariance matrices of hidden-unit activities [for gradient descent ($\gamma \rightarrow 0, \eta = 0$)] exhibit progressive differentiation in the internal representation of hierarchical semantic structure.

- 1 - Thanks to R1 for the improved notation and wording. We agree O'Reilly's work is highly relevant and we should have
2 (and will) cite/discuss Leabra. We will also link to the repository containing code for replicating all results.
- 3 - R2 had two technical concerns: (1) For progressive differentiation (paper Fig. 2a), R2 suggested an alternative measure
4 normalized by overall learning speed. We had investigated several alternatives like this before settling on our metric.
5 Consider Rebuttal Fig. 1a. These settings both show qualitatively similar progressive differentiation, which is captured
6 by our mean lag metric. Further dividing by the time to learn the smallest mode would make the $\gamma < 0$ case appear
7 much less stage-like, because the overall learning rate is faster for $\gamma > 0$ but slower for $\gamma < 0$ (notably it is not simply
8 determined by the absolute value of γ). We found this normalization would make the 2D-map visualization unintuitive.
9 We will clarify these points in the paper. (2) R2 argued that there is no learning when $\gamma = 0$: First, we note that the
10 target is still clamped at the output for $\gamma = 0$ (see blue nodes in paper Fig. 1a) and could participate in learning. We
11 will clarify in paper. Second, however, it is well known that CHL [Eq. (3)] is undefined when $\gamma = 0$ due to division by
12 γ . When we wrote $\gamma = 0$, we always meant $\gamma \rightarrow 0$, i.e., the limit where top-down feedback is infinitesimal. We will
13 correct these labels to be more precise. This limit $\gamma \rightarrow 0, \eta = 0$ reduces to standard gradient descent in which target
14 info is back-propagated to enable learning, not that there is no target info. We agree with R2 and will tone down the
15 bio-plausibility of the CHL network, and our work has no intention to compare the bio-plausibility of different CHL per
16 se (Detorakis et al., 1999). R2 pointed out the similarity of our work to Saxe et al. (2019): Our work is indeed indebted to
17 Saxe et al., but addresses a key limitation of their study. How distinctive are the phenomena that they identify and
18 associate with gradient descent? We show that in fact a swath of learning rules give qualitatively similar behavior.
- 19 - R3 argued our 2D space is too restrictive: In fact we believe a major contribution of our work is proposing a minimalist
20 space that nevertheless encompasses five commonly discussed learning rules. Our metrics do provide a methodology
21 that can be used to characterize any desired learning rule. We will include the references on three-factor learning
22 suggested by R3 since they are indeed relevant. R3 questioned the significance of paper Fig. 6. Fig. 6 shows the degree
23 to which hidden-layer features are learned via error backprop vs. unsupervised Hebbian learning, a key distinction in
24 many theories of neural learning. γ and η do scale the update size, but this does not directly translate to the integrated
25 changes. E.g., if error is driven near zero, the CHL component will stop learning even with large γ . The integrated
26 synaptic strength changes in the network model correspond to important training-induced plasticity that can be measured
27 in electrophysiological experiments (Ahissar & Hochstein, 2004). R3 also argued that illusory correlations depend on
28 the inversion symmetry. We explain here why this is not correct: The input-output map is invariant to this symmetry and
29 our metric directly measures the illusory correlation in the output layer, and hence is not dependent on the SVD. Also
30 we note that nonlinear models (whose dynamics cannot be described by the SVD) show similar phenomena (Rogers &
31 McClelland, 2003).
- 32 - R2 and R3, deep network terminology: Saxe et al. showed that for deep linear networks, the real difference is between
33 0 vs. 1 hidden layer, and we are following this terminology. We will clarify in the text. R4: We do not claim that
34 Hebbian learning can never yield progressive differentiation. We agree that in other tasks (particularly those where the
35 unsupervised statistics are hierarchical too), progressive differentiation could occur. Our goal here was to start with a
36 task environment well studied in prior work, that specifically does not have the target structure embedded in the inputs.
- 37 - R1 and R4 commented on the generalizability of our results to more realistic data sets and tasks beyond the ones we
38 used in the paper. We have now simulated learning dynamics from pixels using distinct images as inputs and have
39 seen similar results (Rebuttal Fig. 1b). The network is trained to learn novel semantic properties (e.g., "cells are made
40 of a molecule X") assigned to perceptually distinct images representing skin samples of hypothetical alien creatures.
41 We use this description of the stimuli because it is a cover story that could fit directly to a human cognitive-learning
42 experiment on this topic (one example of how our theory could inspire or guide future experiments). Multivariate
43 pattern analysis could reveal the evolution of the similarity structure of neural population vectors (recorded using
44 fMRI/MEG) in response to different stimuli over the course of semantic learning. These neural data could be compared
45 to the various theoretical predictions from our framework.