

Deep Predictive Learning as a Model of Human Learning

Randall C. O'Reilly^{a,1}, Jacob L. Russin^a, and John Rohrlich^a

^aDepartment of Psychology, Computer Science, and Center for Neuroscience, University of California Davis

1 **How does the human brain learn new concepts from raw sensory experience, without explicit instruction? This longstanding mystery remains**
2 **unsolved, despite recent demonstrations of the impressive learning power of deep convolutional neural networks (DCNN's), which notoriously**
3 **require explicit training from massive human-labeled datasets. The plausibility of the error backpropagation powering these models has also**
4 **long been questioned on biological grounds, although various related biologically plausible mechanisms have been proposed. Here, we**
5 **show that a biologically based form of *predictive* error-driven learning, where error signals arise from differences between a prediction**
6 **and what actually occurs, learns to systematically categorize 3D objects according to invariant shape properties from raw visual inputs**
7 **alone. We found that these categories match human judgments on the same stimuli, and are consistent with neural representations in**
8 **inferotemporal (IT) cortex in primates. Biologically, we propose that distinctive patterns of connectivity between the neocortex and thalamus**
9 **drive alternating top-down prediction and bottom-up outcome representations over the pulvinar nucleus, at the alpha frequency (10 Hz),**
10 **with the temporal difference driving error-driven learning throughout neocortex. We show that comparison predictive DCNN models lacking**
11 **these biological features did not learn object categories that go beyond the visual input structure. Thus, we argue that incorporating these**
12 **biological properties of the brain can potentially provide a better understanding of human learning at multiple levels relative to existing DCCN**
13 **models.**

Computational Modeling | Predictive Learning | Object Recognition | Pulvinar | Neocortex

1 **T**he fundamental epistemological conundrum of how knowledge emerges from raw experience has
2 plagued philosophers and scientists for centuries. Computational models with powerful learning
3 mechanisms driven by raw images or other sensory inputs provide an attractive way to approach this
4 problem, yet many of the current models based on deep convolutional neural networks (DCNN's)
5 notoriously require explicit training from massive human-labeled datasets (1–3). Such models are
6 cognitively implausible, as non-human primates and human infants learn to recognize and categorize
7 objects without the benefit of such labeled data (4). Furthermore, the biological plausibility of the
8 core learning mechanism, *error backpropagation* (5), has also long been questioned on biological
9 grounds (6), although various related biologically plausible mechanisms have been proposed (7–9).

10 Here we propose a form of *predictive* error-driven learning (10, 11) that learns directly on raw
11 sensory inputs without the need for explicit human-generated labels. This learning mechanism
12 leverages distinctive patterns of connectivity between the neocortex and thalamus (12) (Figure 1) to
13 achieve a biologically based form of predictive learning. In contrast to existing predictive learning
14 frameworks (13–16), we suggest that error signals, as differences between a prediction and what
15 actually occurs, remain as a *temporal difference* in activation states in the network, and are not
16 explicitly represented through error-coding neurons. Specifically, the pulvinar nucleus of the thalamus
17 receives both top-down predictions and bottom-up sensory outcome signals, alternating within an
18 *alpha* frequency cycle (10 Hz, 100 msec), via two distinctive pathways. Thus, our framework has
19 many testable differences from these existing theories, and we argue that existing data is more
20 consistent with our framework.

21 Through large-scale simulations based on the known structure of the visual system, we found that
22 this biologically based predictive learning mechanism developed high-level abstract representations
23 that systematically categorize 3D objects according to invariant shape properties, based on raw visual
24 inputs alone. We found that these categories match human judgments on the same stimuli, and are
25 consistent with neural representations in inferotemporal (IT) cortex in primates (17). Furthermore,
26 we show that comparison predictive DCNN models lacking these biological features (18) did not
27 learn object categories that go beyond the visual input structure. Thus, we argue that incorporating

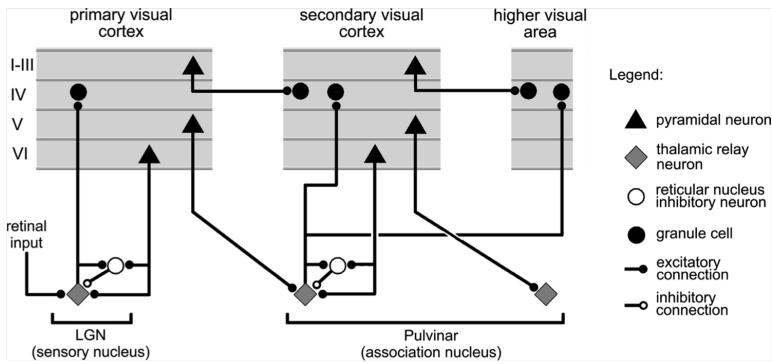


Fig. 1. Summary figure from Sherman & Guillery (2006) showing the strong feedforward driver projection emanating from layer 5IB cells in lower layers (e.g., V1), and the much more numerous feedback “modulatory” projection from layer 6CT cells. We interpret these same connections as providing a prediction (6CT) vs. outcome (5IB) activity pattern over the pulvinar.

these biological properties of the brain can potentially provide a better understanding of human learning at multiple levels relative to existing DCCN models.

Figure 1 shows the thalamocortical circuits characterized by Sherman & Guillery (12) and others, which have two distinct projections converging on the principal thalamic relay cells (TRCs) of the pulvinar (which is interconnected with all higher-level posterior cortical visual areas; (19)). The numerous, weaker projections originating in deep layer VI of the neocortex (the 6CT corticothalamic projecting cells) appear ideal for establishing a top-down prediction state in the pulvinar, based on extensive learning in this pathway and the deep cortical layers that drive it. In contrast, the very sparse (typically one-to-one; (20, 21)) and very strong *driver* inputs originate from lower-level layer V intrinsic bursting cells (5IB), and these can provide a *phasic*, strong bottom-up *ground truth* signal against which the top-down prediction is compared. The 5IB neurons burst at the alpha frequency (22–24), providing a natural timing to the overall predictive learning cycle, consistent with the large and growing literature on alpha properties and effects on perception (25–28).

Based on this and other biological evidence, we hypothesize that this distinctive thalamocortical circuit supports predictive error-driven learning in a way that shapes learning throughout the posterior neocortex (29) (Figure 2a). Specifically, sensory predictions in posterior neocortex are generated roughly every 100 msec at the alpha rhythm, and the pulvinar represents this top-down

Significance Statement

We present a significant advance in understanding how the human brain learns, based on the idea that canonical circuits between the neocortex and thalamus drive alternating phases of prediction and bottom-up outcomes, and the resulting prediction errors (as differences in activation states over time) can drive powerful learning. Critically, we show for the first time that learning based solely on predicting raw visual inputs can generate higher-level abstract categorical representations of 3D objects, which previously has required explicit human-labeled training. This captures the seemingly magic way in which human learning can create knowledge out of raw experience, without explicit teaching.

RCO developed the model, performed the non-PredNet simulations, and drafted the paper. JLR performed the PredNet simulations and analysis, and edited the paper. JR contributed to developing the model and edited the paper.

R. C. O'Reilly is Chief Scientist at eCortex, Inc., which may derive indirect benefit from the work presented here.

¹To whom correspondence should be addressed. E-mail: oreilly@ucdavis.edu

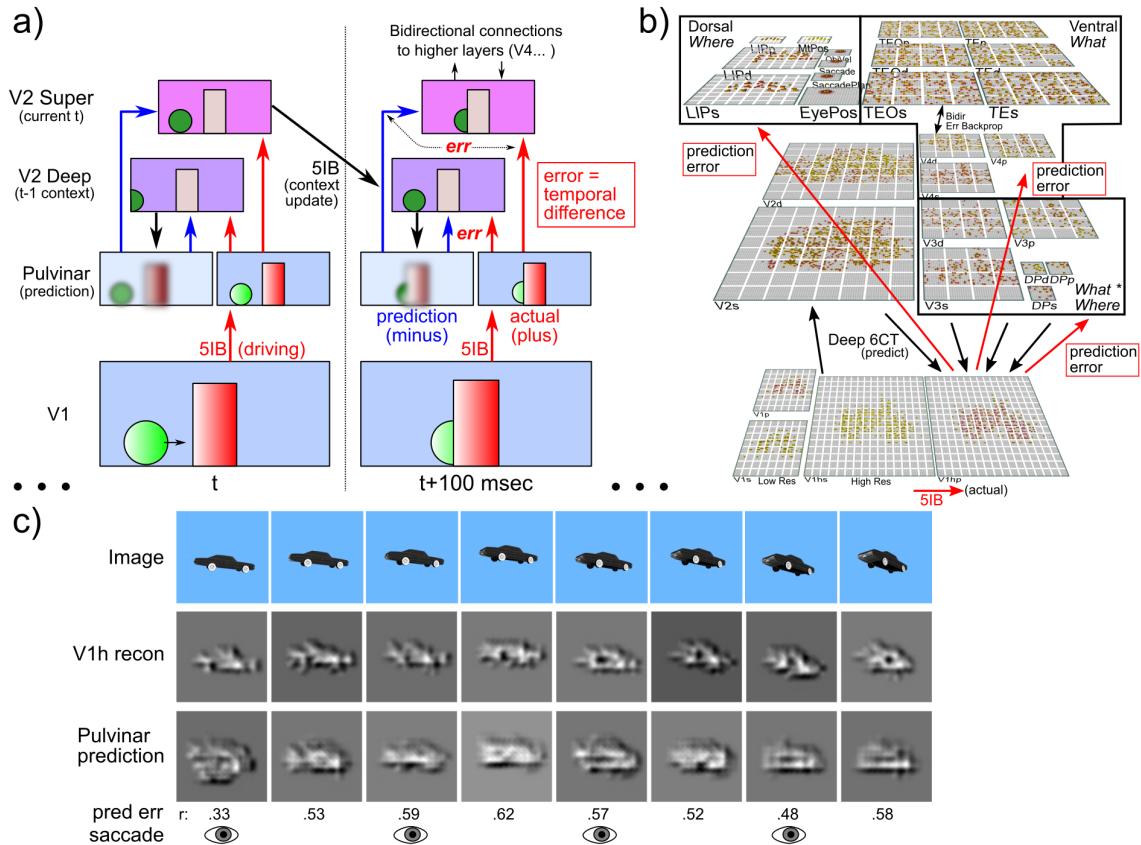


Fig. 2. **a)** Temporal evolution of information flow in the DeepLeabra algorithm predicting visual sequences, over two alpha cycles of 100 msec each. In each alpha cycle, the V2 Deep layer (lamina 5, 6) uses the prior 100 msec of context to generate a prediction (*minus* phase) on the pulvinar thalamic relay cells (TRC). The bottom-up outcome is driven by V1 5IB strong driver inputs (*plus* phase); error-driven learning occurs as a function of the *temporal difference* between these phases, in both superficial (lamina 2, 3) and deep layers, sent via broad pulvinar projections. 5IB bursting in V2 drives update of temporal context in V2 Deep layers, and also the *plus* phase in higher area TRC, to drive higher-level predictive learning. See supporting information (SI) for more details. **b)** The *What-Where-Integration*, *WWI* model. The dorsal *Where* pathway learns first, using easily-abstracted *spatial blobs*, to predict object location based on prior motion, visual motion, and saccade efferent copy signals. This drives strong top-down inputs to lower areas with accurate spatial predictions, leaving the *residual* error concentrated on *What* and *What * Where* integration. The V3 and DP (dorsal prelunate) constitute the *What * Where* integration pathway, binding features and locations. V4, TEO, and TE are the *What* pathway, learning abstracted object category representations, which also drive strong top-down inputs to lower areas. *s* suffix = superficial, *d* = deep, *p* = pulvinar. **c)** Example sequence of 8 alpha cycles that the model learned to predict, with the reconstruction of each image based on the V1 gabor filters (V1 recon), and model-generated prediction (correlation r prediction error shown). The low resolution and reconstruction distortion impair visual assessment, but r values are well above the r 's for each V1 state compared to the previous time step (mean = .38, min of .16 on frame 4 – see SI for more analysis). Eye icons indicate when a saccade occurred.

45 prediction for roughly 75 msec of the alpha cycle as it develops, after which point the layer 5IB
46 intrinsic-bursting neurons send strong, bottom-up driving input to the pulvinar, representing the
47 actual sensory stimulus. Critically, the prediction error is implicit in the temporal difference
48 between these two periods of activity within the alpha cycle over the pulvinar, which is consistent
49 with the biologically plausible form of error-driven cortical learning used in our models (7). The
50 pulvinar sends broad projections back up to all of the areas that drive top-down predictions into
51 it (19, 30), thus broadcasting this error signal to drive local synaptic plasticity in the neocortex.
52 This mathematically approximates gradient descent to minimize overall prediction errors (7). This
53 computational framework makes sense of otherwise puzzling anatomical and physiological properties
54 of the cortical and thalamic networks (12), and is consistent with a wide range of detailed neural
55 and behavioral data (29).

56 A critical question for predictive learning is whether it can develop high-level, abstract ways of
57 representing the raw sensory inputs, while learning from nothing but predicting these low-level visual
58 inputs. For instance, can predictive learning really eliminate the need for human-labeled image
59 datasets where abstract category information is explicitly used to train object recognition models
60 via error-backpropagation? Existing predictive-learning models based on error backpropagation (18)
61 have not demonstrated the development of abstract, categorical representations. Previous work has
62 shown that predictive learning can be a useful method for pretraining networks that are subsequently
63 trained using human-generated labels, but here we focus on the formation of systematic categories
64 *de-novo*.

65 To determine if our biologically based predictive learning model (Figure 2b) can naturally form
66 such categorical encodings in the complete absence of external category labels, we showed the model
67 brief movies of 156 3D object exemplars drawn from 20 different basic-level categories (e.g., car,
68 stapler, table lamp, traffic cone, etc.) selected from the CU3D-100 dataset (31). The objects moved
69 and rotated in 3D space over 8 movie frames, where each frame was sampled at the alpha frequency
70 (Figure 2c). There were also saccadic eye movements every other frame, introducing an additional
71 predictive-learning challenge. An efferent copy signal enabled full prediction of the effects of the
72 eye movement, and allows the model to capture *predictive remapping* (a widely-studied signature of
73 predictive learning in the brain) (32, 33), and introduces additional predictive-learning challenge.
74 The only learning signal available to the model was a prediction error generated by the temporal
75 difference between what it predicted to see in the next frame and what was actually seen.

76 We performed a representational similarity analysis (RSA) on the learned activity patterns at each
77 layer in the model, and found that the highest IT layer (TE) produced a systematic organization of
78 the 156 3D objects into 5 categories (Figure 3a), which visually correspond to the overall shape of
79 the objects (pyramid-shaped, vertically-elongated, round, boxy / square, and horizontally-elongated).
80 This organization of the objects matches that produced by humans making shape similarity judgments
81 on the same set of objects, using the V1 reconstruction as shown in Figure 2c to capture the model's
82 coarse-grained perception (Figure 3b; see supporting information for methods and further analysis).
83 Critically, Figure 3c shows that the overall similarity structure present in IT layers (TEO, TE) of the
84 biological model is significantly different from the similarity structure at the level of the V1 primary
85 visual input. Thus the model, despite being trained only to generate accurate visual input-level
86 predictions, has learned to represent these objects in an abstract way that goes beyond the raw
87 input-level information. Furthermore, this abstract category organization reflects the overall visual
88 shapes of the objects as judged by human participants, suggesting that the model is extracting
89 geometrical shape information that is invariant to the differences in motion, rotation, and scaling
90 that are present in the V1 visual inputs. We further verified that at the highest IT levels in the

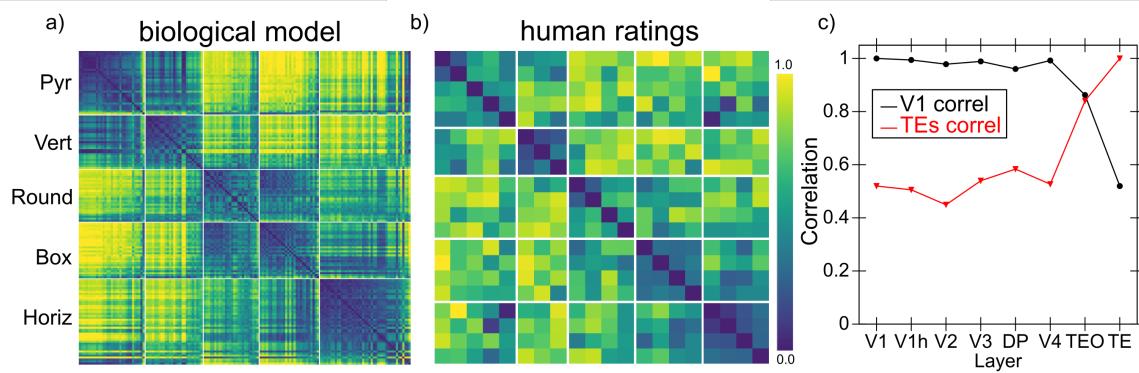


Fig. 3. a) Category similarity structure that developed in the highest layer, TE, of the biologically based predictive learning model, showing 1-correlation similarity of the TE representation for each 3D object against every other 3D object (156 total objects). Blue cells have high similarity, and model has learned block-diagonal clusters or categories of high-similarity groupings, contrasted against dissimilar off-diagonal other categories. Clustering maximized average *within - between* correlation distance (see SI). All items from the same basic-level object categories ($N=20$) are reliably subsumed within learned categories. b) Human similarity ratings for the same 3D objects, presented with the V1 reconstruction (see Fig 1c) to capture coarse perception in model, aggregated by 20 basic-level categories. Each cell is 1 - proportion of time given object pair was rated more similar than another pair (see SI). The human matrix shares the same centroid categorical structure as the model (confirmed by permutation testing and agglomerative cluster analysis, see SI). c) Emergence of abstract category structure over the hierarchy of layers. Red line = correlation similarity between the TE similarity matrix (shown in panel a) and all layers; black line shows correlation similarity between V1 against all layers (1 = identical; 0 = orthogonal). Both show that IT layers (TEO, TE) progressively differentiate from raw input similarity structure present in V1, and, critically, that the model has learned structure beyond that present in the input.

model, a consistent, spatially-invariant representation is present across different views of the same object (e.g., the average correlation across frames within an object was .901). This is also evident in Figure 3a by virtue of the close similarity across multiple objects within the same category.

Further evidence for the progressive nature of representation development in our model is shown in Figure 4, which compares the similarity structures in layers V4 and IT in macaque monkeys (17) with those in corresponding layers in our model. In both the monkeys and our model, the higher IT layer builds upon and clarifies the noisier structure that is emerging in the earlier V4 layer. Considerable other work has also compared DCNN representations with these same data from monkeys (17), but it is essential to appreciate that those DCNN models were explicitly trained on the category labels, making it somewhat less than surprising that such categorical representations developed. By contrast, we reiterate that our model has discovered its categorical representations entirely on its own, with no explicit categorical inputs or training of any kind.

Figure 5 shows the results from a purely backpropagation-based (Bp) version of the same model architecture, and a standard PredNet model (18) with extensive hyperparameter optimization (see SI). In the Bp model, the highest layers in the network form a simple binary category structure overall, and the detailed item-level similarity structure does not diverge significantly from that present at the lowest V1 inputs, indicating that it has not formed novel systematic structured representations, in contrast to those formed in the biologically based model. Similar results were found in the PredNet model, where the highest layer representations remained very close to the V1 input structure. Thus, it is clear that the additional biologically derived properties are playing a critical role in the development of abstract categorical representations that go beyond the raw visual inputs. These properties include: excitatory bidirectional connections, inhibitory competition, and an additional Hebbian form of learning that serves as a regularizer (similar to weight decay) on top of predictive error-driven learning (34, 35).

Each of these properties could promote the formation of categorical representations. Bidirectional connections enable top-down signals to consistently shape lower-level representations, creating significant attractor dynamics that cause the entire network to settle into discrete categorical attractor states. By contrast, backpropagation networks typically lack these kinds of attractor dynamics, and this could contribute significantly to their relative lack of categorical learning.

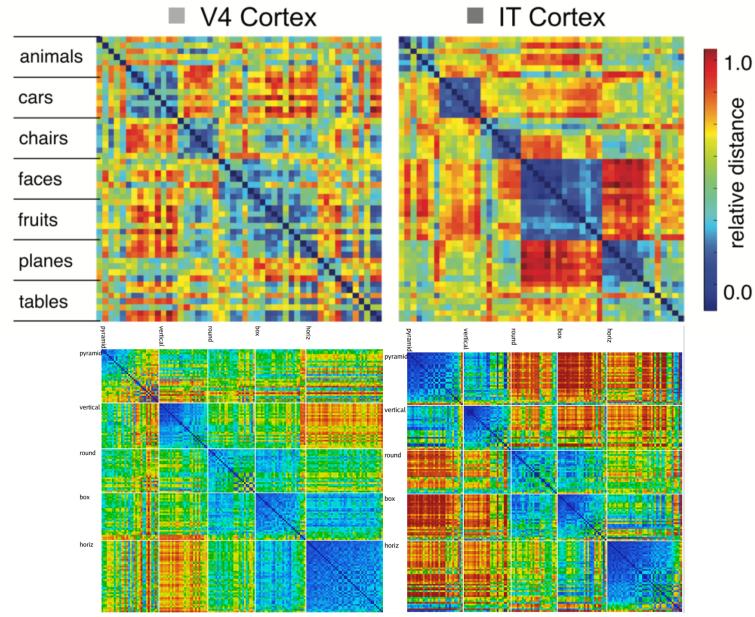


Fig. 4. Comparison of progression from V4 to IT in macaque monkey visual cortex (top row, from Cadieu et al., 2014) versus same progression in model (replotted using comparable color scale). Although the underlying categories are different, and the monkeys have a much richer multi-modal experience of the world to reinforce categories such as foods and faces, the model nevertheless shows a similar qualitative progression of stronger categorical structure in IT, where the block-diagonal highly similar representations are more consistent across categories, and the off-diagonal differences are stronger and more consistent as well (i.e., categories are also more clearly differentiated). Note that the critical difference in our model versus those compared in Cadieu et al. 2014 and related papers is that they explicitly trained their models on category labels, whereas our model is *entirely self-organizing* and has no external categorical training signal.

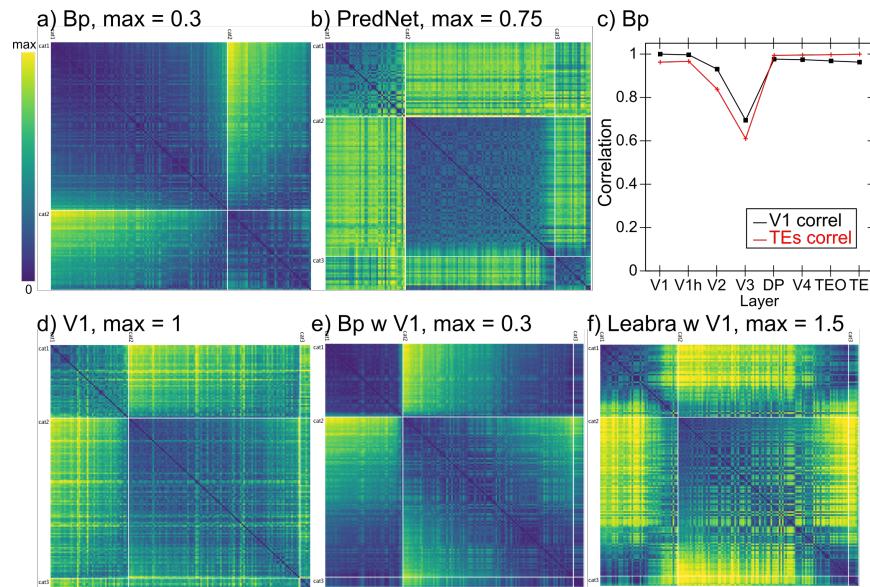


Fig. 5. a) Best-fitting category similarity for TE layer of the backpropagation (Bp) model with the same What / Where structure as the biological model. Only two broad categories are evident, and the lower *max* distance (0.3 vs. 1.5 in biological model) means that the patterns are highly similar overall. b) Best-fitting similarity structure for the PredNet model, in the highest of its layers (layer 6), which is more differentiated than Bp ($\text{max} = 0.75$) but also less cleanly similar within categories (i.e., less solidly blue along the block diagonal), and overall follows a broad category structure similar to V1. c) Comparison of similarity structures across layers in the Bp model (compare to Figure 2c): unlike in the biological model, the V1 structure is largely preserved across layers, and is little different from the structure that best fits the TE layer shown in panel a, indicating that the model has not developed abstractions beyond the structure present in the visual input. Layer V3 is most directly influenced by spatial prediction errors, so it differs from both in strongly encoding position information. d) The best fitting V1 structure, which has 2 broad categories and banana is in a third category by itself. The lack of dark blue on the block diagonal indicates that these categories are relatively weak, and every item is fairly dissimilar from every other. e) The same similarities shown in panel a for Bp TE also fit reasonably well sorted according to the V1 structure (and they have a similar average within - between contrast differences, of 0.0838 and 0.0513 – see SI for details). f) The similarity structure from the biological model resorted in the V1 structure does *not* fit well: the blue is not aligned along the block diagonal, and the yellow is not strictly off-diagonal. This is consistent with the large difference in average contrast distance: 0.5071 for the best categories vs. 0.3070 for the V1 categories.

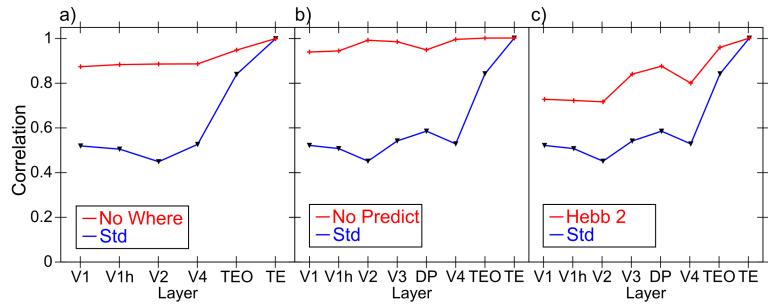


Fig. 6. Effects of various manipulations on the extent to which TE representations differentiate from V1. *Std* is the same result shown in Figure 2c from the intact model for ease of comparison. All of the following manipulations significantly impair the development of abstract TE categorical representations (i.e., TE is more similar V1 and the other layers). **a)** Dorsal *Where* pathway lesions, including lateral inferior parietal sulcus (LIP), V3, and dorsal prelunate (DP). This pathway is essential for regressing out location-based prediction errors, so that the residual errors concentrate feature-encoding errors that train the *What* pathway. **b)** Allowing the deep layers full access to current-time information, thus effectively eliminating the prediction demand and turning the network into an auto-encoder, which significantly impairs representation development, and supports the importance of the challenge of predictive learning for developing deeper, more abstract representations. **c)** Reducing the strength of Hebbian learning by 20% (from 2.5 to 2), demonstrating the essential role played by this form of learning on shaping categorical representations. Eliminating Hebbian learning entirely (not shown) prevented the model from learning anything at all, as it also plays a critical regularization and shaping role on learning.

120 Hebbian learning drives the formation of representations that encode the principal components
 121 of activity correlations over time, which can help more categorical representations coalesce (and
 122 results below already indicate its importance). Inhibition, especially in combination with Hebbian
 123 learning, drives representations to specialize on more specific subsets of the space. Ongoing
 124 work is attempting to determine which of these is essential in this case (perhaps all of them) by
 125 systematically introducing some of these properties into the backpropagation model, though this is
 126 difficult because full bidirectional recurrent activity propagation, which is essential for conveying
 127 error signals top-down in the biological network, is incompatible with the standard efficient form of
 128 error backpropagation, and requires much more computationally intensive and unstable forms of fully
 129 recurrent backpropagation (36, 37). Furthermore, Hebbian learning requires inhibitory competition
 130 which is difficult to incorporate within the backpropagation framework.

131 Figure 6 shows just a few of the large number of parameter manipulations that have been conducted
 132 to develop and test the final architecture. For example, we hypothesized that separating the overall
 133 prediction problem between a spatial *Where* vs. non-spatial *What* pathway (38, 39), would strongly
 134 benefit the formation of more abstract, categorical object representations in the *What* pathway.
 135 Specifically, the *Where* pathway can learn relatively quickly to predict the overall spatial trajectory
 136 of the object (and anticipate the effects of saccades), and thus effectively regress out that component
 137 of the overall prediction error, leaving the residual error concentrated in object feature information,
 138 which can train the ventral *What* pathway to develop abstract visual categories. Figure 6a shows that,
 139 indeed, when the *Where* pathway is lesioned, the formation of abstract categorical representations
 140 in the intact *What* pathway is significantly impaired. Figure 6b shows that full predictive learning,
 141 as compared to just encoding and decoding the current state (which is much easier computationally,
 142 and leads to much better overall accuracy), is also critical for the formation of abstract categorical
 143 representations — prediction is a “desirable difficulty” (40). Finally, Figure 6c shows the impact of
 144 reducing Hebbian learning, which impairs category learning as expected.

145 In conclusion, we have demonstrated that learning based strictly on predicting what will be seen
 146 next is, in conjunction with a number of critical biologically motivated network properties and
 147 mechanisms, capable of generating abstract, invariant categorical representations of the overall
 148 shapes of objects. The nature of these shape representations closely matches human shape similarity
 149 judgments on the same objects. Thus, predictive learning has the potential to go beyond the surface
 150 structure of its inputs, and develop systematic, abstract encodings of the “deeper” structure of the

environment. Relative to existing machine-learning-based approaches in “deep learning”, which have generally focused on raw categorization accuracy measures using explicit category labels or other human-labeled inputs, the results here suggest that focusing more on the nature of what is learned in the model might provide a valuable alternative approach. Considerable evidence in cognitive neuroscience suggests that the primary function of the many nested (“deep”) layers of neural processing in the neocortex is to *simplify* and aggressively *discard* information (41), to produce precisely the kinds of extremely valuable abstractions such as object categories, and, ultimately, symbol-like representations that support high-level cognitive processes such as reasoning and problem-solving (42, 43). Thus, particularly in the domain of predictive or generative learning, the metric of interest should not be the accuracy of prediction itself (which is indeed notably worse in our biologically based model compared to the DCNN-based PredNet and backpropagation models), but rather whether this learning process results in the formation of simpler, abstract representations of the world that can in turn support higher levels of cognitive function.

Considerable further work remains to be done to more precisely characterize the essential properties of our biologically motivated model necessary to produce this abstract form of learning, and to further explore the full scope of predictive learning across different domains. We strongly suspect that extensive cross-modal predictive learning in real-world environments, including between sensory and motor systems, is a significant factor in infant development and could greatly multiply the opportunities for the formation of higher-order abstract representations that more compactly and systematically capture the structure of the world (44). Future versions of these models could thus potentially provide novel insights into the fundamental question of how deep an understanding a pre-verbal human, or a non-verbal primate, can develop (11, 45), based on predictive learning mechanisms. This would then represent the foundation upon which language and cultural learning builds, to shape the full extent of human intelligence.

ACKNOWLEDGMENTS. We thank Dean Wyatte, Tom Hazy, Seth Herd, Kai Krueger, Tim Curran, David Sheinberg, Lew Harvey, Jessica Mollick, Will Chapman, Helene Devillez, and the rest of the CCN Lab for many helpful comments and suggestions. Supported by: ONR grants ONR N00014-19-1-2684 / N00014-18-1-2116, N00014-14-1-0670 / N00014-16-1-2128, N00014-18-C-2067, N00014-13-1-0067, D00014-12-C-0638. This work utilized the Janus supercomputer, which is supported by the National Science Foundation (award number CNS-0821794) and the University of Colorado Boulder. The Janus supercomputer is a joint effort of the University of Colorado Boulder, the University of Colorado Denver and the National Center for Atmospheric Research. All data and materials will be available at <https://github.com/ccnlab/deep-obj-cat> upon publication.

1. Krizhevsky A, Sutskever I, Hinton GE (2012) ImageNet Classification with Deep Convolutional Neural Networks in *Advances in Neural Information Processing Systems 25*, eds. Pereira F, Burges CJC, Bottou L, Weinberger KQ. (Curran Associates, Inc.), pp. 1097–1105.
2. LeCun Y, Bengio Y, Hinton G (2015) Deep learning. *Nature* 521(7553):436–444.
3. Schmidhuber J (2015) Deep learning in neural networks: An overview. *Neural Networks* 61:85–117.
4. Lake BM, Ullman TD, Tenenbaum JB, Gershman SJ (2017) Building machines that learn and think like people. *Behavioral and Brain Sciences* 40.
5. Rumelhart DE, Hinton GE, Williams RJ (1986) Learning representations by back-propagating errors. *Nature* 323(9):533–536.
6. Crick F (1989) The recent excitement about neural networks. *Nature* 337:129–132.
7. O'Reilly RC (1996) Biologically plausible error-driven learning using local activation differences: The generalized recirculation algorithm. *Neural Computation* 8(5):895–938.
8. Xie X, Seung HS (2003) Equivalence of backpropagation and Contrastive Hebbian Learning in a layered network. *Neural Computation* 15(2):441–454.
9. Bengio Y, Mesnard T, Fischer A, Zhang S, Wu Y (2017) STDP-compatible approximation of backpropagation in an energy-based model. *Neural Computation* 29(3):555–577.
10. Elman JL (1990) Finding Structure In Time. *Cognitive Science* 14(2):179–211.
11. Elman J, et al. (1996) *Rethinking Innateness: A Connectionist Perspective on Development*. (MIT Press, Cambridge, MA).
12. Sherman S, Guillery R (2006) *Exploring the Thalamus and Its Role in Cortical Function*. (MIT Press, Cambridge, MA).
13. Mumford D (1992) On the computational architecture of the neocortex. II. The role of cortico-cortical loops. *Biological Cybernetics* 66(3):241–251.
14. Rao RP, Ballard DH (1999) Predictive coding in the visual cortex: A functional interpretation of some extra-classical receptive-field effects. *Nature Neuroscience* 2(1):79–87.
15. Kawato M, Hayakawa H, Inui T (1993) A forward-inverse optics model of reciprocal connections between visual cortical areas. *Network: Computation in Neural Systems* 4(4):415–422.
16. Friston K (2005) A theory of cortical responses. *Philosophical Transactions of the Royal Society B* 360(1456):815–836.
17. Cadieu CF, et al. (2014) Deep Neural Networks Rival the Representation of Primate IT Cortex for Core Visual Object Recognition. *PLoS Computational Biology* 10(12):e1003963.
18. Lotter W, Kreiman G, Cox D (2016) Deep predictive coding networks for video prediction and unsupervised learning. *arXiv:1605.08104 [cs, q-bio]*.
19. Shipp S (2003) The functional logic of cortico-pulvinar connections. *Philosophical Transactions of the Royal Society of London B* 358(1438):1605–1624.
20. Rockland KS (1998) Convergence and branching patterns of round, type 2 corticopulvinar axons. *The Journal of Comparative Neurology* 390(4):515–536.
21. Rockland KS (1996) Two types of corticopulvinar terminations: Round (type 2) and elongate (type 1). *The Journal of comparative neurology* 368:57–87.

- 206 22. Lorincz ML, Kekesi KA, Juhasz G, Crunelli V, Hughes SW (2009) Temporal framing of thalamic relay-mode firing by phasic inhibition during the alpha rhythm. *Neuron* 63(5):683–696.
- 207 23. Franceschetti S, et al. (1995) Ionic mechanisms underlying burst firing in pyramidal neurons: Intracellular study in rat sensorimotor cortex. *Brain Research* 696(1–2):127–139.
- 208 24. Saalmann YB, Pinsk MA, Wang L, Li X, Kastner S (2012) The pulvinar regulates information transmission between cortical areas based on attention demands. *Science* 337(6095):753–756.
- 209 25. Buffalo EA, Fries P, Landman R, Buschman TJ, Desimone R (2011) Laminar differences in gamma and alpha coherence in the ventral stream. *Proceedings of the National Academy of Sciences of the United States of America* 108(27):11262–11267.
- 211 26. VanRullen R, Koch C (2003) Is perception discrete or continuous? *Trends in Cognitive Sciences* 7(5):207–213.
- 212 27. Jensen O, Bonnefond M, VanRullen R (2012) An oscillatory mechanism for prioritizing salient unattended stimuli. *Trends in Cognitive Sciences* 16(4):200–206.
- 213 28. Fiebelkorn IC, Kastner S (2019) A rhythmic theory of attention. *Trends in Cognitive Sciences* 23(2):87–101.
- 214 29. O'Reilly RC, Wyatte D, Rohrlich J (2014) Learning Through Time in the Thalamocortical Loops. *arXiv:1407.3432 [q-bio]*.
- 215 30. Mumford D (1991) On the computational architecture of the neocortex. *Biological Cybernetics* 65(2):135–145.
- 216 31. O'Reilly RC, Wyatte D, Herd S, Mingus B, Jilk DJ (2013) Recurrent Processing during Object Recognition. *Frontiers in Psychology* 4(124).
- 217 32. Duhamel JR, Colby CL, Goldberg ME (1992) The updating of the representation of visual space in parietal cortex by intended eye movements. *Science* 255(5040):90–92.
- 218 33. Cavanagh P, Hunt AR, Afraz A, Rolfs M (2010) Visual stability based on remapping of attention pointers. *Trends in Cognitive Sciences* 14(4):147–153.
- 219 34. O'Reilly RC (1998) Six Principles for Biologically-Based Computational Models of Cortical Cognition. *Trends in Cognitive Sciences* 2(11):455–462.
- 220 35. O'Reilly RC, Munakata Y (2000) *Computational Explorations in Cognitive Neuroscience: Understanding the Mind by Simulating the Brain*. (MIT Press, Cambridge, MA).
- 221 36. Williams RJ, Zipser D (1992) Gradient-based learning algorithms for recurrent networks and their computational complexity in *Backpropagation: Theory, Architectures and Applications*, eds. Chauvin Y, Rumelhart DE. (Erlbaum, Hillsdale, NJ).
- 222 37. Pineda FJ (1987) Generalization of Backpropagation to Recurrent Neural Networks. *Physical Review Letters* 18:2229–2232.
- 224 38. Ungerleider LG, Mishkin M (1982) Two Cortical Visual Systems in *The Analysis of Visual Behavior*, eds. Ingle DJ, Goodale MA, Mansfield RJW. (MIT Press, Cambridge, MA), pp. 549–586.
- 225 39. Goodale MA, Milner AD (1992) Separate visual pathways for perception and action. *Trends in Neurosciences* 15(1):20–25.
- 226 40. Bjork RA (1994) Memory and metamemory considerations in the training of human beings in *Metacognition: Knowing about Knowing*. (The MIT Press, Cambridge, MA, US), pp. 185–205.
- 227 41. Simons DJ, Rensink RA (2005) Change blindness: Past, present, and future. *Trends in cognitive sciences* 9(1):16–20.
- 228 42. Rougier NP, Noelle D, Braver TS, Cohen JD, O'Reilly RC (2005) Prefrontal Cortex and the Flexibility of Cognitive Control: Rules Without Symbols. *Proceedings of the National Academy of Sciences* 102(20):7338–7343.
- 230 43. O'Reilly RC, et al. (2014) How Limited Systematicity Emerges: A Computational Cognitive Neuroscience Approach in *The Architecture of Cognition: Rethinking Fodor and Pylyshyn's Systematicity Challenge*, eds. Calvo IP, Symons J. (MIT Press, Cambridge, MA).
- 231 44. Yu C, Smith LB (2012) Embodied attention and word learning by toddlers. *Cognition* 125(2):244–262.
- 232 45. Spelke E, Breinlinger K, Macomber J, Jacobson K (1992) Origins of Knowledge. *Psychological Review* 99(4):605–632.

Supporting Information For: Deep Predictive Learning as a Model of Human Learning

Randall C. O'Reilly, Jacob L. Russin, and John Rohrlich

Correspondence to: oreilly@ucdavis.edu

November 8, 2019

This PDF includes:

Materials and Methods

Figures S1-S9

Table S1

All of the materials described here, including the experimental study, the computational models, and the code to perform the representational similarity analysis, are all available on our github account at: <https://github.com/ccnlab/deep-obj-cat> For the computational models in particular, the most complete understanding can only be had by directly examining the code for the models, as there are a number of details that are not efficiently captured in this supplementary materials text.

1 Representational Similarity Analysis Methods

The different representations being compared here are:

Leabra: The DeepLeabra (biological model) TE layer representations (specifically TEs = superficial – results are very similar for deep as well).

Bp: The TEs layer representations from the backpropagation version of biological model, including *What*, *Where* and *What * Where* integration layers, trained with the V1p and V1hp (low and high resolution pulvinar) layers as final output layers, using the time t target pattern from the $t - 1$ input (i.e., as a predictive network).

V1: The gabor-filtered representation of the visual input to both of the above models, which was identical across them.

PredNet: Highest layer (6th Layer) of the PredNet architecture.

Expt: Similarity matrix constructed from human pairwise similarity judgments (see *Behavioral Experiment Methods*).

An optimal category cluster can be defined as one that has high within-cluster similarity and low between-cluster similarity. This can be operationalized by the *contrast* distance metric, based on a 1-correlation (*correlation distance*) measure, as the difference between the average within-cluster similarity and the average between-cluster similarity:

$$cd = \langle 1 - r_{in} \rangle - \langle 1 - r_{out} \rangle \quad (1)$$

With distance-like 1-correlation values, this contrast distance should be minimized (it is typically negative), or equivalently the contrast on raw correlation values can be maximized (it is typically a positive number – just the sign flip of distance value). We refer to the positive numbers and maximization here as that is more natural.

Starting with an initial set of clusters, a permutation-based hill-climbing strategy was used to determine a local minimum in this measure: each item was tested in each of the other possible categories, and if that configuration reduced the overall average contrast distance metric across all items, then it was adopted and the process iterated until no such permutation improved the metric. This algorithm can only decrease the number of clusters (by moving all items out of a given cluster), so different numbers of initial clusters can be used to search the overall space.

Figure S1 shows the resulting categories. The Bp model converged on the same cluster state from all starting configurations tested, varying from 5 to 2 initial categories. This is the cluster set

Centroid	Bp		
1. pyramid	3. round cont'd	1. cat1	1. cat1 cont'd
• banana	• handgun	• banana	• handgun
• layercake	• chair	• layercake	• chair
• trafficcone	4. box	• trafficcone	• slrcamera
• sailboat	• slrcamera	• sailboat	• elephant
• trex	• elephant	• trex	• piano
2. vertical	• piano	• person	• fish
• person	• fish	• guitar	• car
• guitar	5. horiz	• tablelamp	2. cat2
• tablelamp	• car	• doorknob	• heavycannon
3. round	• heavycannon	• donut	• stapler
• doorknob	• stapler		• motorcycle
• donut	• motorcycle		
V1		PredNet	
1. cat1	2. cat2 cont'd	1. cat1	2. cat2 cont'd
• trafficcone	• handgun	• trafficcone	• slrcamera
• sailboat	• slrcamera	• sailboat	• elephant
• person	• elephant	• person	• fish
• guitar	• piano	• guitar	• car
• tablelamp	• fish	• tablelamp	• heavycannon
• chair	• car	• layercake	• stapler
2. cat2	• heavycannon	2. cat2	• motorcycle
• layercake	• stapler		
• trex	• motorcycle		
• doorknob	3. cat3		
• donut	• banana		

Figure S1: Shape categories used for similarity matrix plots in main paper. *Centroid* shape categories are near-best for both the Leabra model and the Expt results, and fit our visual intuitions about overall shape. *Bp* are reliably optimal for Bp model from all starting points. *V1* are reliably optimal for V1 inputs, and also were close to the best for the Bp and PredNet layer 6 representations. *PredNet* are best stable solution for PredNet layer 6.

shown in Figure 5a of the main paper, and has an average contrast distance (*acd*) of 0.0838 (this is relatively low because the patterns were overall quite similar). Likewise, the V1 patterns (which were the same across Leabra and Bp models) reliably converged on the same pattern (shown in Figure 5d), with *acd* = 0.2448.

For the PredNet layer 6 representations, starting from the V1 categories gave the best results of any other set ($acd = 0.1967$), and a few permutations resulted in a reliable solution that was arrived at from all other 3 category starting points tested, shown in Figure S1 ($acd = 0.2820$). This indicates that PredNet did not go much beyond the structure present in the input, even though it did not use the V1 gabor filtering used in the Leabra and Bp models (i.e., this V1-level encoding well-captures the structure of the visual inputs in general). The PredNet pixel and layer 1 representations both converged on essentially a single monolithic category with very low acd (0.0018, 0.0013).

For the Leabra TE representations, we found a set of *centroid* shape categories that are near-best when considering both the Leabra model and the results from the human behavioral experiment (Expt). Starting from these categories, the permutation analysis converged on reducing the size of the vertical and round categories to one item each, over a sequence of 5 steps. This is consistent with the observation from Figure 3a that there are three broader categories within which the 5 finer-grained categories are embedded (i.e., vertical and pyramid are overall similar to each other, as are round and box). Nevertheless, our initial visual intuition about the broad shape categories, along with a bias against having single-item categories, reinforced the use of the finer-grained centroid selection. The average contrast difference of our centroid selection is 0.5071, while the maximal result from the permutation was 0.5526, which is a relatively small proportional difference.

Furthermore, once we had collected the human experimental data (*Expt*), it was clear that it strongly coincided with our original shape intuitions, and with the finer-grained 5 category centroid structure. Starting from the centroid categories, the maximal permutation made only 3 changes, moving trex (T-rex) and handgun into the horizontal category, and chair into the pyramid, going from a distance score of 0.3083 to 0.3225, which is a relatively small improvement. However, using the maximal *Expt* clusters directly on the Leabra model gives a lower acd measure of 0.3745 (compared to 0.5071 for centroid), so the centroid categories represent a good middle-ground between experiment and the model, and this strong shared similarity structure with near-optimal cluster structures confirms that the model and people are encoding largely the same information.

In contrast, if we organize the experiment similarity matrix using the Bp categories, it produces a very poor average contrast distance measure of 0.0643 (compared to 0.3083 for the centroid categories), strongly suggesting that people's shape representations are not compatible with that simple structure.

Another approach to determining clusters from similarity matrices, *agglomerative clustering*, starts with all items as singletons, and iteratively combines the closest two into a new cluster. The results for the Leabra and Expt similarity matrices are shown in Figure S2, which has also color-coded the items in terms of their category status according to the centroid structure. Due to a strong history dependency in the clustering process, and the indeterminacy of reducing a high-dimensional similarity structure down to two dimensions, structure beyond the leaf level is not very reliable (ties are also broken by a random number generator), but nevertheless you can clearly see that in both cases items from the same cluster are almost always together as leaves in the plots. This then provides additional converging support for the idea that the model is learning the same kind of shape categories as people have.

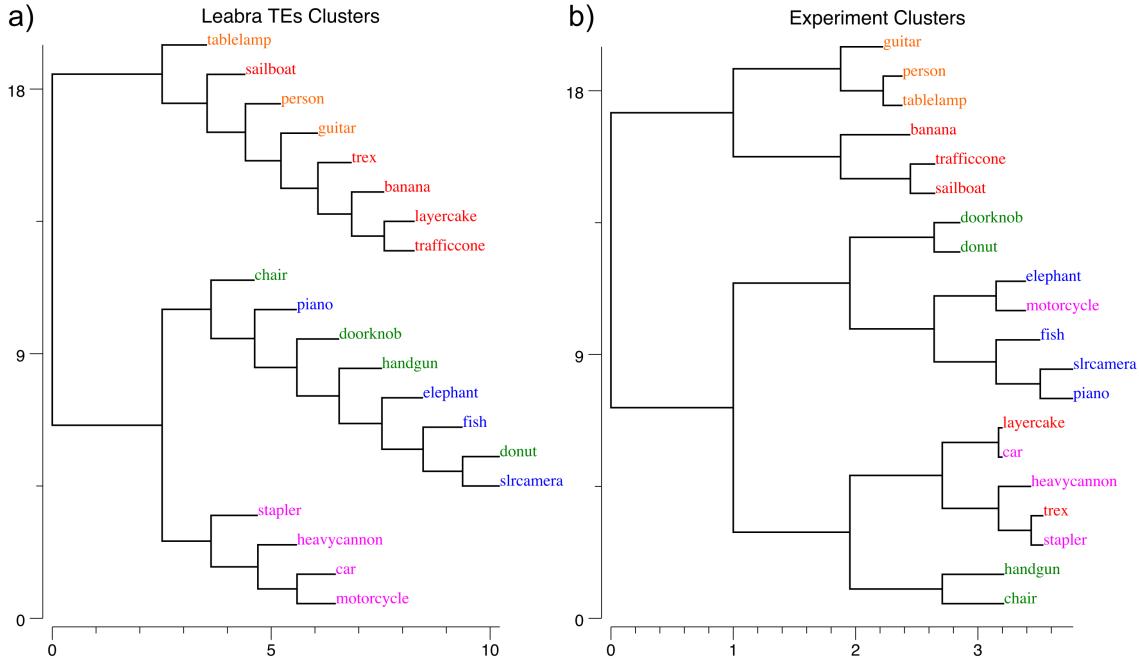


Figure S2: Agglomerative clustering on the Leabra and Expt representations, with the centroid categories color coded. The most reliable information from this is the leaf-level groupings, as the rest of the structure is indeterminant and history dependent in reducing higher-dimensional structure down to a 2D plot. Both cluster plots show a strong tendency to group leaf items together in the same centroid categories, with a few exceptions in each case. Also, the Leabra plot nicely captures the broader 3-category structure evident in the similarity matrix plots, within which the 5 finer-grained centroid categories are organized. Overall, this provides further confirmation that the model and the human subjects are organizing the shapes in largely the same way.



Figure S3: Example stimulus from the behavioral experiment, using the V1 reconstruction of the actual input images presented to the model, to better capture the coarse-grained perception of the model. Subjects were requested to choose which of the two pairs, Left or Right, was most similar in terms of *overall shape*.

2 Behavioral Experiment Methods

The behavioral experiment was conducted on Amazon.com's MTurk web platform under University of Colorado IRB approval (19-0176), using 30 participants each categorizing up to 800 image pairs as shown in Figure S3, using the standard *simple image categorization* framework with a lightly customized script. Objects were drawn from the 156 3D object set, but data was aggregated

in terms of the 20 basic-level categories (car, stapler, etc) because we could not sample all 156 x 156 object pairs. Thus, the resulting data was aggregated for each category pair in terms of the proportion of times when that pair was selected when presented.

The individual images were produced by reconstructing from the V1 transform that the computational model used in its high resolution V1 input layer, to give human participants as similar of an experience as possible to how the model “saw” the objects, and to reduce the influence of existing semantic knowledge which was entirely missing in our model (Figure S3).

3 Biological Model Methods

This section provides more information about the *DeepLeabra What-Where Integration (WWI)* model. The purpose of this information is to give more detailed insight into the model’s function beyond the level provided in the main text, but with a model of this complexity, the only way to really understand it is to explore the model itself. It is available for download at: <https://github.com/ccnlab/deep-obj-cat/sims/C++> Furthermore, the best way to understand this model is to understand the framework in which it is implemented, which is explained in great detail, with many running simulations explaining specific elements of functionality, at <http://ccnbook.colorado.edu>

3.1 Layer Sizes and Structure

Figure 2 in the main text shows the general configuration of the model, and Table S1 shows the specific sizes of each of the layers, and where they receive inputs from.

All the activation and general learning parameters in the model are at their standard Leabra defaults.

3.2 Projections

The general principles and patterns of connectivity are shown in Figure S4 (and Figure 1 in the main text).

Detailed each of the specific parameters associated with the different projections shown in Table S1 would take too much space — those interested in this level of detail should download the model from the link shown above. There are topographic projections between many of the lower-level retinotopically-mapped layers, consistent with our earlier vision models [1]. For example the 8x8 unit groups in V2 are reduced down to the 4x4 groups in V3 via a 4x4 unit-group topographic projection, where neighboring units have half-overlapping receptive fields (i.e., the field moves over 2 unit groups in V2 for every 1 unit group in V3), and the full space is uniformly tiled by using a wrap-around effect at the edges. Similar patterns of connectivity are used in current deep convolutional neural networks. However, we do not share weights across units as in a true convolutional network.

The projections from ObjVel (object velocity) and SaccadePlan layers to LIPs, LIPd were initialized with a topographic sigmoidal pattern that moved as a function of the position of the unit group, by a factor of .5, while the projections from EyePos were initialized with a gaussian pattern.

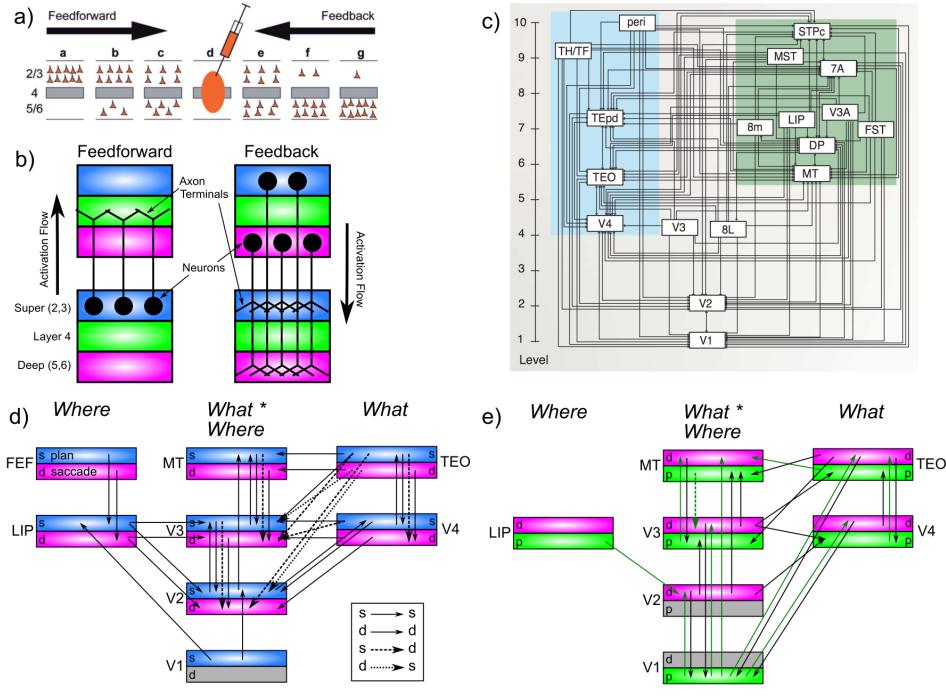


Figure S4: Principles of connectivity in DeepLeabra. **a)** Markov et al (2014) data showing density of *retrograde* labeling from a given injection in a middle-level area (d): most feedforward projections originate from superficial layers of lower areas (a,b,c) and deep layers predominantly contribute to feedback (and more strongly for longer-range feedback). **b)** Summary diagram showing most feedforward connections originating in superficial layers of lower area, and terminating in layer 4 of higher area, while feedback connections can originate in either superficial or deep layers, and in both cases terminate in both superficial and deep layers of the lower area (adapted from Felleman & Van Essen, 1991). **c)** Anatomical hierarchy as determined by percentage of superficial layer source labeling (SLN) by Markov et al (2014) — the hierarchical levels are well matched for our model, but we functionally divide the dorsal pathway (shown in green background) into the two separable components of a *Where* and a *What * Where* integration pathway. **d)** Superficial and deep-layer connectivity in the model. Note the repeating motif between hierarchically-adjacent areas, with bidirectional connectivity between superficial layers, and feedback into deep layers from both higher-level superficial and deep layers, according to canonical pattern shown in panels a and b. Special patterns of connectivity from TEO to V3 and V2, involving crossed super-to-deep and deep-to-super pathways, provide top-down support for predictions based on high-level object representations. **e)** Connectivity for deep layers and pulvinar in the model, which generally mirror the corticocortical pathways (in d). Each pulvinar layer (p) receives 5IB driving inputs from the labeled layer (e.g., V1p receives 5IB drivers from V1). In reality these neurons are more distributed throughout the pulvinar, but it is computationally convenient to organize them together as shown. Deep layers (d) provide predictive input into pulvinar, and pulvinar projections send error signals (via temporal differences between predictions and actual state) to *both* deep and superficial layers of given areas (only d shown). Most areas send deep-layer prediction inputs into the main V1p prediction layer, and receive reciprocal error signals therefrom. The strongest constraint we found was that pulvinar outputs (colored green) must generally project only to higher areas, not to lower areas, with the exceptions of DPp → V3 and LIPp → V2. V2p was omitted because it is largely redundant with V1p in this simple model.

Area	Name	Units		Pools		Receiving Projections
		X	Y	X	Y	
V1	V1s	4	5	8	8	
	V1p	4	5	8	8	V1s V2d V3d V4d TEOd
V1h	V1hs	4	5	16	16	
	V1hp	4	5	16	16	V1s V2d V3d V4d TEOd
Eyes	EyePos	21	21			
	SaccadePlan	11	11			
	Saccade	11	11			
Obj	ObjVel	11	11			
V2	V2s	10	10	8	8	V1s LIPs V3s V4s TEOd V1p V1hp
	V2d	10	10	8	8	V2s V1p V1hp LIPd LIPp V3d V4d V3s TEOs
LIP	MtPos	1	1	8	8	V1s
	LIPs	4	4	8	8	MtPos ObjVel SaccadePlan EyePos LIPp
	LIPd	4	4	8	8	LIPs LIPp ObjVel Saccade EyePos
	LIPp	1	1	8	8	MtPos V1s LIPd
V3	V3s	10	10	4	4	V2s V4s TEOs DPs LIPs V1p V1hp DPp TEOd
	V3d	10	10	4	4	V3s V1p V1hp DPp LIPd DPd V4d V4s DPs TEOs
	V3p	10	10	4	4	V3s V2d DPd TEOd
DP	DPs	10	10			V2s V3s TEOs V1p V1hp V3p TEOp
	DPd	10	10			DPs V1p V1hp DPp TEOd
	DPp	10	10			DPs V2d V3d DPd TEOd
V4	V4s	10	10	4	4	V2s TEOs V1p V1hp
	V4d	10	10	4	4	V4s V1p V1hp V4p TEOd TEOs
	V4p	10	10	4	4	V4s V2d V3d V4d TEOd
TEO	TEOs	10	10	4	4	V4s V1p V1hp TEs
	TEOd	10	10	4	4	TEOs TEOd V1p V1hp V4p TEOp TEp TED
	TEOp	10	10	4	4	TEOs V3d V4d TEOd TED
TE	TEs	10	10	4	4	TEOs V1p V1hp
	TED	10	10	4	4	TEs TED V1p V1hp V4p TEOp TEp TEOd
	TEp	10	10	4	4	TEs V3d V4d TEOd

Table S1: Layer sizes, showing numbers of units in one pool (or entire layer if Pool is missing), and the number of Pools of such units, along X,Y axes. Each area has three associated layers: *s* = superficial layer, *d* = deep layer (context updated by 51B neurons in same area, shown in bold), *p* = pulvinar layer (driven by 51B neurons from associated area, shown in bold).

These patterns multiplied uniformly distributed random weights in the .25 to .75 range, with the lowest values in the topographic pattern having a multiplier of .6, while the highest had a multiplier of 1 (i.e., a fairly subtle effect). This produced faster convergence of the LIP layer when doing *Where* pathway pre-training compared to purely random initial weights. In addition to exploring different patterns of overall connectivity, we also explored differences in the relative strengths of receiving projections, which can be set with a `wt_scale.rel` parameter in the simulator. All feedforward pathways have a default strength of 1. For the feedback projections, which are typically weaker (consistent with the biology), we explored a discrete range of strengths, typically

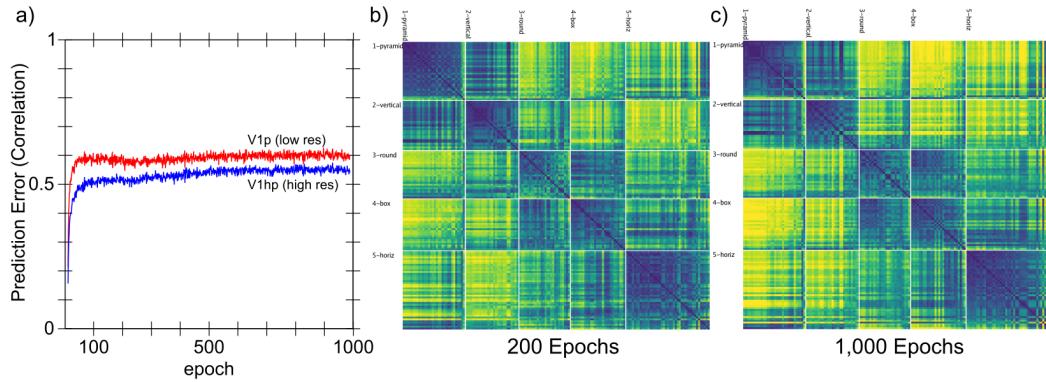


Figure S5: **a)** Predictive learning curve for DeepLeabra, showing the correlation between prediction and actual over the two different V1 layers. Initial learning is quite rapid, followed by a slower but progressive learning process that reflects development of the IT representations (e.g., manipulations that interfere with those areas selectively impair this part of the learning curve). Overall prediction accuracy remains far from perfect, as shown in Figure 2c in main text, and significantly worse than the backpropagation-based models. This is a typical finding from Leabra models which are significantly more constrained as a result of bidirectional attractor dynamics, Hebbian learning, and inhibitory competition – i.e., the very things that are likely important for forming abstract categorical representations. **b)** Similarity matrix over TEs layer at 200 epochs, which has less contrast and definition compared to the 1,000 epoch result (**c** also shown in Figure 3a in main text).

.5, .2, .1, and .05. The strongest top-down projections were into V2s from LIP and V3, while most others were .2 or .1. Likewise projections from the pulvinar were weaker, typically .1. These differences in strength sometimes had large effects on performance during the initial bootstrapping of the overall model structure, but in the final model they are typically not very consequential for any individual projection.

3.3 Training Parameters

Training typically consisted of 512 alpha trials per epoch (51.2 seconds of real time equivalent), for 1,000 such epochs. Each trial was generated from a virtual reality environment in the emergent simulator, that rendered first-person views with moving eye position onto the object tumbling through space with fixed motion and rotation parameters over the sequence of 8 frames (see Figure 2c in main text for representative example). Because the start of each sequence of 8 frames is unpredictable, we turned off learning for that trial, which improves learning overall. We have recently developed an automatic such mechanism based on the running-average (and running variance) of the prediction error, where we turn off learning whenever the current prediction error z-normalized by these running average values is below 1.5 standard deviations, which works well, and will be incorporated into future models. Biologically, this could correspond to a connection between pulvinar and neuromodulatory areas that could regulate the effective learning rate in this way.

Figure S5a shows the learning trajectory of the model, indicating that it learns quite rapidly. This rapid initial learning is likely facilitated by the extensive use of shortcut connections covering

from all over the simulated visual system onto the V1 pulvinar layers, and direct projections back from these pulvinar layers. Thus, error signals are directly communicated and can drive learning quickly and efficiently. However, there are also extensive indirect, bidirectional connections among the superficial layers, which can drive indirect error backpropagation learning as well.

3.4 Testing Parameters

To be able to monitor similarity metrics as the model trained, we used a running-average integration of neural activity across trials to accumulate the patterns used for the RSA analysis described above. Specifically, for each object, and each frame, the current activation pattern across each layer was recorded and averaged unit-by-unit with a time constant of $\tau = 10$. Critically, by integrating separately for each frame, this running-average computation did not introduce any bias for temporally-adjacent frames to be more similar. Nevertheless, when we computed the frame-to-frame similarities for TE, they were quite high (.901 correlation on average across all objects).

3.5 Model Algorithms

The biologically-based model was implemented using the Leabra framework, which is described in detail in previous publications [2, 3, 4, 5], and summarized here. There are two main implementations of Leabra, one in the C++ *emergent* software, and a new one using Go and Python language at: <https://github.com/emer/leabra>. There are also other simpler implementations in Python and MATLAB, see <https://grey.colorado.edu/emergent/index.php/Leabra>. Both of the preceding links contain a full detailed description of the algorithm. These same equations and standard parameters have been used to simulate over 40 different models in [2, 3], and a number of other research models. Thus, the model can be viewed as an instantiation of a systematic modeling framework using standardized mechanisms, instead of constructing new mechanisms for each model. Here, we only detail properties of the predictive learning algorithm that go beyond the basic Leabra model.

3.5.1 Deep Context

At the end of every plus phase, a new deep-layer context net input is computed from the dot product of the context weights times the sending activations, just as in the standard net input:

$$\eta = \langle x_i w_{ij} \rangle = \frac{1}{n} \sum_i x_i w_{ij} \quad (2)$$

This net input is then added in with the standard net input at each cycle of processing.

The relative strength of these context layer inputs was set progressively larger for higher layers in the network, with a maximum of 4 in V4, TEO, and TE. In addition, TEO and TE received *self* context projections which provide an extended window of temporal context into the prior 200 msec interval. These self projections were connected only within the narrower Pool level of units, enabling these neurons to develop mutually-excitatory loops to sustain activations over the multiple trials when the same object was present. We hypothesize that these modifications correspond to

biological adaptations in IT cortex that likewise support greater sustained activation of object-level representations.

Learning of the context weights occurs as normal, but using the sending activation states from the *prior* time step's activation.

3.5.2 Computational and Biological Details of SRN-like Functionality

Predictive auto-encoder learning has been explored in various frameworks, but the most relevant to our model comes from the application of the SRN to a range of predictive learning domains [6, 7]. One of the most powerful features of the SRN is that it enables error-driven learning, instead of arbitrary parameter settings, to determine how prior information is integrated with new information. Thus, SRNs can learn to hold onto some important information for a relatively long interval, while rapidly updating other information that is only relevant for a shorter duration. This same flexibility is present in our DeepLeabra model. Furthermore, because this temporal context information is hypothesized to be present in the deep layers throughout the entire neocortex (in every microcolumn of tissue), the DeepLeabra model provides a more pervasive and interconnected form of temporal integration compared to the SRN, which typically just has a single temporal context layer associated with the internal “hidden” layer of processing units.

An extensive computational analysis of what makes the SRN work as well as it does, and explorations of a range of possible alternative frameworks, has led us to an important general principle: *subsequent outcomes determine what is relevant from the past*. At some level, this may seem obvious, but it has significant implications for predictive learning mechanisms based on temporal context. It means that the information encoded in a temporal context representation cannot be learned at the time when that information is presently active. Instead, the relevant contextual information is learned on the basis of what happens next. This explains the peculiar power of the otherwise strange property of the SRN: the temporal context information is preserved as a *direct copy* of the state of the hidden layer units on the previous time step (Figure S6), and then learned synaptic weights integrate that copied context information into the next hidden state (which is then copied to the context again, and so on). This enables the error-driven learning taking place in the *current* time step to determine how context information from the *previous* time step is integrated. And the simple direct copy operation eschews any attempt to shape this temporal context itself, instead relying on the learning pressure that shapes the hidden layer representations to also shape the context representations. In other words, this copy operation is essential, because there is no other viable source of learning signals to shape the nature of the context representation itself (because these learning signals require future outcomes, which are by definition only available later).

The direct copy operation of the SRN is however seemingly problematic from a biological perspective: how could neurons copy activations from another set of neurons at some discrete point in time, and then hold onto those copied values for a duration of 100 msec, which is a reasonably long period of time in neural terms (e.g., a rapidly firing cortical neuron fires at around 100 Hz, meaning that it will fire 10 times within that context frame). However, there is an important transformation of the SRN context computation, which is more biologically plausible, and compatible with the structure of the deep network (Figure S6). Specifically, instead of copying an entire set of activation states, the context activations (generated by the phasic 5IB burst) are immediately sent through the adaptive synaptic weights that integrate this information, which we think occurs in the

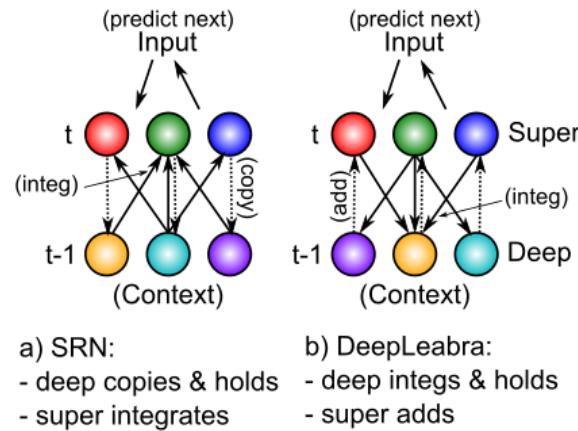


Figure S6: How the DeepLeabra temporal context computation compares to the SRN mathematically. **a)** In a standard SRN, the context (deep layer biologically) is a copy of the hidden activations from the prior time step, and these are held constant while the hidden layer (superficial) units integrate the context through learned synaptic weights. **b)** In DeepLeabra, the deep layer performs the weighted integration of the soon-to-be context information from the superficial layer, and then holds this integrated value, and feeds it back as an additive net-input like signal to the superficial layer. The context net input is pre-computed, instead of having to compute this same value over and over again. This is more efficient, and more compatible with the diffuse interconnections among the deep layer neurons. Layer 6 projections to the thalamus and back recirculate this pre-computed net input value into the superficial layers (via layer 4), and back into itself to support maintenance of the held value.

6CC (corticocortical) and other lateral integrative connections from 5IB neurons into the rest of the deep network. The result is a *pre-computed net input* from the context onto a given hidden unit (in the original SRN terminology), not the raw context information itself. Computationally, and metabolically, this is a much more efficient mechanism, because the context is, by definition, unchanging over the 100 msec alpha cycle, and thus it makes more sense to pre-compute the synaptic integration, rather than repeatedly re-computing this same synaptic integration over and over again (in the original feedforward backpropagation-based SRN model, this issue did not arise because a single step of activation updating took place for each context update — whereas in our bidirectional model many activation update steps must take place per context update).

There are a couple of remaining challenges for this transformation of the SRN. First, the pre-computed net input from the context must somehow persist over the subsequent 100 msec period of the alpha cycle. We hypothesize that this can occur via NMDA and mGluR channels that can easily produce sustained excitatory currents over this time frame. Furthermore, the reciprocal excitatory connectivity from 6CT to TRC and back to 6CT could help to sustain the initial temporal context signal. Second, these contextual integration synapses require a different form of learning algorithm that uses the sending activation from the prior 100 msec, which is well within the time constants in the relevant calcium and second messenger pathways involved in synaptic plasticity.

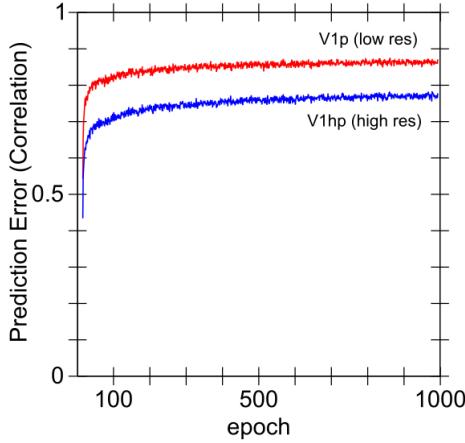


Figure S7: Learning curves for the backpropagation version of the WWI model. Although it achieves better predictive accuracy than the DeepLeabra version, it fails to acquire abstract object category structure, indicating a potential tradeoff between simplifying and categorizing inputs, versus predicting precisely where the low-level visual features will move.

4 Backpropagation Model Methods

The backpropagation version of the WWI model has exactly the same layer sizes and *feedforward* patterns of connectivity as the DeepLeabra version. Topographically, the V1p and V1hp pulvinar layers serve as output layers at the highest level of the network, receiving all the various connections from deep layers as shown in Table S1. Likewise, the LIPp served as a target output layer for the Where pathway. To achieve predictive learning, the V1 pulvinar targets were from the scene at time t , while the V1s inputs were from the scene at time $t - 1$. We also ran a comparison auto-encoder model that had inputs and target outputs from the same time step, and it showed even less systematic organization of its higher-level representations, further supporting the notion that predictive learning is important, across all frameworks. The learning curve for the predictive version is shown in Figure S7, which shows better overall prediction accuracy compared to the DeepLeabra model. However, as the RSA showed, this backpropagation model failed to learn object categories that go beyond the input similarity structure, indicating that perhaps it was paying too much “attention” in learning to this low-level structure, and lacked the necessary mechanisms to enable it to impose a simplifying higher-level structure on top of these inputs.

5 PredNet Model Methods

The PredNet architecture was designed to incorporate principles from predictive coding theory into a neural network model for predicting the next frame in a video sequence. Details of the model can be found in the original paper [8], but here we provide a brief overview of the architecture.

5.1 Architecture

PredNet is a deep convolutional neural network that is composed of layers containing discrete modules. The lowest layer generates a prediction of incoming inputs (i.e. the pixels in the next frame), while each of the higher layers attempts to predict the *errors* made by the previous layer. Each layer contains an input convolutional module (A_l), a recurrent representational module (R_l), a prediction module (\hat{A}_l), and a representation of its own errors (E_l). The input convolutional module (A_l) transforms its input with a set of standard convolutional filters, a rectified linear activation function, and a max-pooling operation. The recurrent representation module (R_l) is a convolutional LSTM, which is a recurrent convolutional network that replaces the matrix multiplications in the standard LSTM equations with convolutions, allowing it to maintain a spatially organized representation of its inputs over time. The prediction module (\hat{A}_l) consists of another standard convolutional layer and rectified linear activation that is used to generate predictions from the output of R_l . These predictions are then compared against the output of the input convolutional module (A_l). The errors generated in this comparison are represented explicitly in E_l , which applies a rectified linear activation to a concatenation of the positive ($A_l - \hat{A}_l$) and negative ($\hat{A}_l - A_l$) prediction errors. These errors then become the inputs to the next layer.

$$A_l^t = \begin{cases} x_t, & \text{if } l = 0 \\ \text{MaxPool}(\text{ReLU}(\text{Conv}(E_{l-1}^t))), & \text{if } l > 0 \end{cases} \quad (3)$$

$$\hat{A}_l^t = \text{ReLU}(\text{Conv}(R_l^t)) \quad (4)$$

$$E_l^t = [\text{ReLU}(A_l^t - \hat{A}_l^t); \text{ReLU}(\hat{A}_l^t - A_l^t)] \quad (5)$$

$$R_l^t = \text{ConvLSTM}(E_l^{t-1}, R_l^{t-1}, \text{UpSample}(R_{l+1}^t)) \quad (6)$$

At each time step in the video sequence, PredNet generates a prediction of the next frame. This is done as follows: first, the R_l is computed for each layer starting from the top of the hierarchy (because each R_l^t depends on input from R_{l+1}^t), and then the A_l^t , \hat{A}_l^t and E_l^t are computed in a feed-forward fashion (because each A_l^t depends on input from the layer below, E_{l-1}^t).

All analyses in the RSA were conducted using the representations from the R_l layers.

5.2 Implementation details

All experiments with the PredNet architecture were performed using PyTorch. An informal hyperparameter search was conducted to find the settings that maximized representational similarity to the human judgments. This was done by conducting RSA on each layer for each hyperparameter setting, and computing, according to the Centroid categories derived from the human data, the difference between the average within-category similarity and the average between-category similarity. Our final architecture had 6 layers with 3, 16, 32, 64, 128, and 256 filters in the A_l and R_l modules, and 3x3 kernels throughout the whole network. We also found that using sigmoid and tanh activation functions in fully-connected convolutional LSTMs slightly improved performance, so these were used for all experiments.

The weights in the PredNet model are trained using error backpropagation. Predictions are generated and errors are computed at all levels of the hierarchy, but the model performs better when

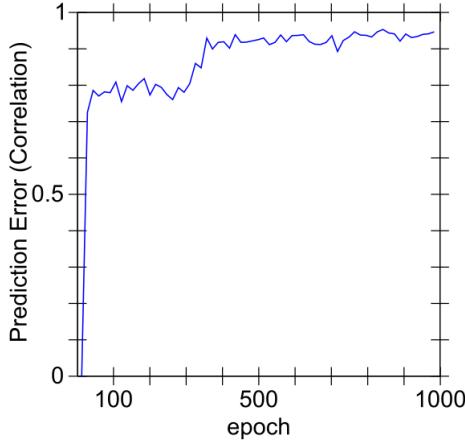


Figure S8: Learning curves for the PredNet model. This model achieves the best overall prediction performance but also has the least well differentiated, categorical representations.

only the lowest layer’s errors are backpropagated [8]. We confirmed these results with experiments that backpropagated the errors in higher layers, in which performance (in terms of mean squared error) was marginally reduced but the RSA results were similar. For this reason, all reported experiments used a PredNet that was trained by only backpropagating the lowest level error.

The model was trained using a batch size of 8 and an Adam optimizer with a learning rate of 0.0001, with no scheduler, for 150,000 batches. A training curve is shown in Figure S8, showing that it achieves the best overall prediction accuracy of any model we tested, and yet does not have representations that are as differentiated or categorical as our biologically based model, as shown in the main paper.

5.3 Regularization experiments

As discussed in the main paper, our biologically based model includes a number of important biologically motivated properties that may be contributing to the development of its categorical representations. These properties, including excitatory bidirectional connections, inhibitory competition, and an additional form of Hebbian learning, may be acting as regularizers that encourage categorical learning. We therefore tested whether standard regularization methods used in deep learning would have similar effects on the representations developed in the PredNet architecture. We tested 1) batch normalization, 2) dropout (0.1, 0.3, and 0.5), and 3) weight decay (0.01, 0.001, 0.0001, 0.00001). All experiments with batch normalization and weight decay showed reduced performance (in terms of both prediction error on the test set and within-category correlation). As shown in figure S9, dropout marginally improved the within-category correlation while also slightly improving prediction accuracy, so a dropout rate of 0.1 was used for the comparison to our biologically based model in the main paper.

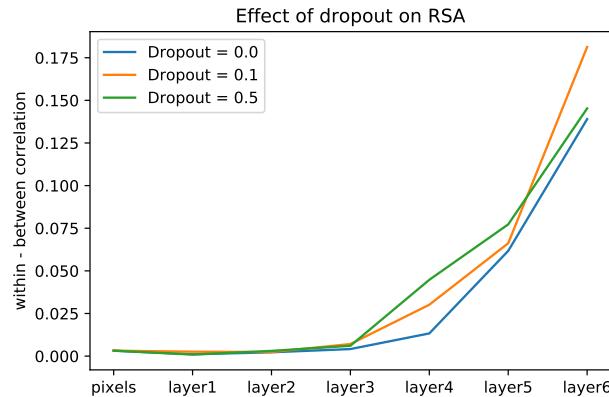


Figure S9: Effect of dropout in PredNet on RSA, as measured by the difference between the average within-category correlation and the average between category correlation (using the Centroid categories derived from human data). Dropout marginally improves the category structure learned in PredNet.

References

- [1] O'Reilly RC, Wyatte D, Herd S, Mingus B, Jilk DJ (2013) Recurrent Processing during Object Recognition. *Frontiers in Psychology* 4(124).
- [2] O'Reilly RC, Munakata Y, Frank MJ, Hazy TE, Contributors (2012) *Computational Cognitive Neuroscience*. (Wiki Book, 1st Edition, URL: <http://ccnbook.colorado.edu>).
- [3] O'Reilly RC, Munakata Y (2000) *Computational Explorations in Cognitive Neuroscience: Understanding the Mind by Simulating the Brain*. (MIT Press, Cambridge, MA).
- [4] O'Reilly RC (1998) Six Principles for Biologically-Based Computational Models of Cortical Cognition. *Trends in Cognitive Sciences* 2(11):455–462.
- [5] O'Reilly RC (1996) Biologically plausible error-driven learning using local activation differences: The generalized recirculation algorithm. *Neural Computation* 8(5):895–938.
- [6] Elman JL (1990) Finding Structure In Time. *Cognitive Science* 14(2):179–211.
- [7] Elman J, et al. (1996) *Rethinking Innateness: A Connectionist Perspective on Development*. (MIT Press, Cambridge, MA).
- [8] Lotter W, Kreiman G, Cox D (2016) Deep predictive coding networks for video prediction and unsupervised learning. *arXiv:1605.08104 [cs, q-bio]*.