



**Citation:** Jiang LP, Rao RPN (2024) Dynamic predictive coding: A model of hierarchical sequence learning and prediction in the neocortex. PLoS Comput Biol 20(2): e1011801. https://doi.org/10.1371/journal.pcbi.1011801

**Editor:** Jonathan Rubin, University of Pittsburgh, UNITED STATES

Received: March 5, 2023

Accepted: January 4, 2024

Published: February 8, 2024

Copyright: © 2024 Jiang, Rao. This is an open access article distributed under the terms of the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

**Data Availability Statement:** Data and code for reproducing all simulations in the paper are available at <a href="https://github.com/lpjiang97/dynamic-predictive-coding.">https://github.com/lpjiang97/dynamic-predictive-coding.</a>

Funding: This work was supported by the National Institutes of Health (1UF1NS126485-01 to RPNR as co-PI), National Science Foundation (NSF) EFRI (2223495 to RPNR as co-PI), Defense Advanced Research Projects Agency (DARPA) Contract (HR001120C0021 to RPNR via a subcontract), a UW + Amazon Science Hub grant, a Weill Neurohub Investigator grant, a Frameworks grant

RESEARCH ARTICLE

# Dynamic predictive coding: A model of hierarchical sequence learning and prediction in the neocortex

Linxing Preston Jiang 1,2,3, Rajesh P. N. Rao 1,2,3 \*

1 Paul G. Allen School of Computer Science & Engineering, University of Washington, Seattle, Washington, United States of America, 2 Center for Neurotechnology, University of Washington, Seattle, Washington, United States of America, 3 Computational Neuroscience Center, University of Washington, Seattle, Washington, United States of America

\* rao@cs.washington.edu

# **Abstract**

We introduce dynamic predictive coding, a hierarchical model of spatiotemporal prediction and sequence learning in the neocortex. The model assumes that higher cortical levels modulate the temporal dynamics of lower levels, correcting their predictions of dynamics using prediction errors. As a result, lower levels form representations that encode sequences at shorter timescales (e.g., a single step) while higher levels form representations that encode sequences at longer timescales (e.g., an entire sequence). We tested this model using a two-level neural network, where the top-down modulation creates low-dimensional combinations of a set of learned temporal dynamics to explain input sequences. When trained on natural videos, the lower-level model neurons developed space-time receptive fields similar to those of simple cells in the primary visual cortex while the higherlevel responses spanned longer timescales, mimicking temporal response hierarchies in the cortex. Additionally, the network's hierarchical sequence representation exhibited both predictive and postdictive effects resembling those observed in visual motion processing in humans (e.g., in the flash-lag illusion). When coupled with an associative memory emulating the role of the hippocampus, the model allowed episodic memories to be stored and retrieved, supporting cue-triggered recall of an input sequence similar to activity recall in the visual cortex. When extended to three hierarchical levels, the model learned progressively more abstract temporal representations along the hierarchy. Taken together, our results suggest that cortical processing and learning of sequences can be interpreted as dynamic predictive coding based on a hierarchical spatiotemporal generative model of the visual world.

# **Author summary**

The brain is adept at predicting stimuli and events at multiple timescales. How do the neuronal networks in the brain achieve this remarkable capability? We propose that the neocortex employs dynamic predictive coding to learn hierarchical spatiotemporal

from the Templeton World Charity Foundation, and a Cherng Jia and Elizabeth Yun Hwang Professorship to RPNR. The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript."

**Competing interests:** The authors have declared that no competing interests exist.

representations. Using computer simulations, we show that when exposed to natural videos, a hierarchical neural network that minimizes prediction errors develops stable and longer timescale responses at the higher level; lower-level neurons learn space-time receptive fields similar to the receptive fields of primary visual cortical cells. The same network also exhibits several effects in visual motion processing and supports cue-triggered activity recall. Our results provide a new framework for understanding the genesis of temporal response hierarchies and activity recall in the neocortex.

### Introduction

The ability to predict future stimuli and event outcomes is critical for perceiving and interacting with a highly dynamic world. At the neural circuit level, predictions could compensate for neural transmission delays and engage with the world in real-time. At the cognitive level, planning a sequence of actions to achieve a desired goal relies on predictions of the sensory consequences of motor commands. These abilities are predicated on two requirements: (a) the brain must infer the dynamics of sensory stimuli to make spatiotemporal predictions based on an internal model of the world, and (b) the brain's temporal representations must span different timescales to support predictions over both short and long horizons.

Many experimental studies have provided evidence for such computations. Predictive representations of upcoming stimuli have been found in various open and closed-loop paradigms where animals developed experience-dependent visual and auditory expectations [1–5]. Other empirical evidence suggests that cortical representations exhibit a hierarchy of timescales and an increase in stability from lower-order to higher-order areas across both sensory and cognitive regions [6–9]. We asked the question: could such phenomena be explained by the neocortex learning a spatiotemporal generative model based on a temporal hierarchy of representations?

Predictive coding provides a unifying framework for understanding perception and prediction in terms of learning hierarchical generative models of the environment [10-14]. Here, we present dynamic predictive coding (DPC), a new predictive coding model for learning hierarchical temporal representations. The central idea of our proposal is that our perceptual system learns temporally abstracted representations that encode entire sequences rather than single points at any given time. Specifically, DPC assumes that higher-level model neurons modulate the transition dynamics of lower-level networks, building on the computational concept of hypernetworks [15]. Hypernetworks are neural networks that generate the parameters (synaptic weights) for another neural network. However, generating an entire set of high-dimensional synaptic weights is not neurally plausible. Instead, DPC models the transition dynamics at a lower level of a hierarchy using a small set of modulation weights for a group of learned transition matrices. These weights implement "top-down" gain modulation of the lower-level synapses [16, 17] and are predicted by the higher level through a feedback network (a hypernetwork) connecting the higher to the lower level. Compared to previous normative models of video processing that either do not learn the temporal dynamics between images [18–22] or presume a fixed temporal hierarchy [23, 24] (see Discussion), the DPC model offers a neural implementation of spatiotemporal prediction that learns the transition dynamics of the input and adapts its hierarchical temporal representation to the intrinsic timescales of the data.

We tested the DPC model using a two-level neural network trained on natural and artificial image sequences to minimize spatiotemporal prediction errors. After training, the lower-level neurons developed space-time receptive fields similar to those found in simple cells in the

primary visual cortex (V1) [25]. Neurons in the second level learned to capture input dynamics on a longer timescale and their responses exhibited greater stability compared to responses in the first level, similar to the temporal response hierarchies observed in the cortex [6-9]. We further show that the learned sequence representations in the network can explain both predictive and postdictive effects seen in visual processing [26-29], reproducing several aspects of the flash-lag illusion [26, 30, 31]. When linked to an associative memory mimicking the role of the hippocampus, the network allowed storage of episodic memories and exhibited cue-triggered activity recall after repeated exposure to a fixed input sequence, an effect previously reported in rodents [1], human V1 [32-34] and monkey V4 [35]. Lastly, when extended to three levels, the top-level neurons learned to encode the transition dynamics of the secondlevel states, which in turn encoded the transition dynamics of the first-level states, thereby yielding a hierarchical temporal representation of input image sequences. Together, these results support the hypothesis that the neocortex uses dynamic predictive coding based on a hierarchical spatiotemporal generative model to learn and interpret input sequences at multiple levels of temporal abstractions. Some of the results presented herein appeared previously in a conference proceedings [36].

#### Results

## Dynamic predictive coding

The DPC model assumes that spatiotemporal inputs are generated by a hierarchical generative model (Fig 1a) (see also [37]). We describe here a two-level hierarchical model (see Discussion for the possibility of extending the model to more levels). The lower level of the model follows the traditional predictive coding model in generating images using a set of spatial filters **U** and a latent state vector  $\mathbf{r}_t$ , which is sparse [38], for each time step t:  $\mathbf{I}_t = \mathbf{U}\mathbf{r}_t + \mathbf{n}$  where  $\mathbf{n}$  is zero mean Gaussian white noise. The temporal dynamics of the state  $\mathbf{r}_t$  is modeled using K learnable transition matrices  $\{\mathbf{V}_k\}_{k=1}^K$  which can be linearly combined using a set of "modulation" weights given by a K-dimensional vector  $\mathbf{w}$ . This vector of weights is generated by the higher-level state vector  $\mathbf{r}^h$  using a function  $\mathcal{H}$  (Fig 1b), implemented as a neural network (a "hypernetwork" [15]—see "Hypernetworks and neural gain modulation" in S1 Text):

$$\mathbf{w} = \mathcal{H}(\mathbf{r}^h) \tag{1}$$

$$\mathbf{V} = \sum_{k=1}^{K} w_k \mathbf{V}_k. \tag{2}$$

Here,  $w_k$  is the kth component of the vector  $\mathbf{w}$ . The lower-level state vector at time t+1 is generated as  $\mathbf{r}_{t+1} = \text{ReLU}(\mathbf{V}\mathbf{r}_t) + \mathbf{m}$  where  $\mathbf{m}$  is zero mean Gaussian white noise. Note that this is one particular parameterization for top-down modulation of the lower-level transition dynamics, with the hypernetwork formulation allowing other types of parameterizations (see "Hypernetworks and neural gain modulation" in S1 Text).

The generative model in Fig 1b can be implemented in a hierarchical neural network: the higher-level state  $\mathbf{r}^h$ , represented by higher-level neurons, generates a top-down modulation  $\mathbf{w}$  via a top-down feedback neural network  $\mathcal{H}$ , and this top-down input  $\mathbf{w}$  influences the groups of lower-level neurons representing  $\mathbf{V}_i$  through gain modulation [16, 17] (see "Hypernetworks and neural gain modulation" in S1 Text for details). We propose that such a computation could be implemented by cortical pyramidal neurons receiving top-down modulation via their apical dendrites (through gain control [17, 39]) and the recurrent state  $\mathbf{r}_t$  (and input prediction errors) via their basal dendrites, and integrating these to predict the next state (Fig 1c).

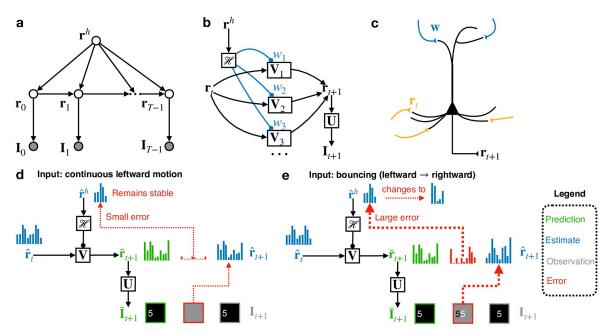


Fig 1. Dynamic predictive coding. (a) Generative model for dynamic predictive coding. (b) Parameterization of the model. The higher-level state modulates the lower-level transition matrices through a top-down network ("hypernetwork")  $\mathcal{H}$ . (c) A possible neural implementation of the generative model using cortical pyramidal neurons. Pyramidal neurons receive the top-down embedding vector input via synapses at apical dendrites and the current recurrent state vector via basal dendrites, and produce as their output the next state vector. (d) Schematic depiction of an inference step when the dynamics at the lower level is stable. The higher-level state remains stable due to minimal prediction errors. (e) Depiction of an inference step when the lower-level dynamics changes. The resulting large prediction errors drive updates to the higher-level state to account for the new lower-level dynamics.

When an input sequence is presented, the model employs a Bayesian filtering approach to perform online inference on the latent vectors [40] by minimizing a loss function that includes prediction errors and penalties from prior distributions over the latent variables (see Methods). Given the model's estimates  $\hat{\mathbf{r}}_t$  and  $\hat{\mathbf{r}}^h$  at time t, the estimate  $\hat{\mathbf{r}}_{t+1}$  of  $\mathbf{r}$  at time t+1 is computed by gradient descent to minimize the sum of the input prediction error  $\|\mathbf{I}_{t+1} - \mathbf{U}\hat{\mathbf{r}}_t\|_2^2$  and the temporal state prediction error  $\|\mathbf{r}_{t+1} - \text{ReLU}(\mathbf{V}\hat{\mathbf{r}}_t)\|_2^2$  plus a sparseness penalty. Similarly, the second level estimates  $\hat{\mathbf{r}}^h$  is updated using the temporal prediction error plus a prior-related penalty. The model's parameters are learned by minimizing the same prediction errors across all time steps and input sequences, further reducing the errors not accounted for by the inference process above for latent vectors (see Methods).

# Hierarchical predictive coding of natural videos

We implemented the DPC model described above using a two-level neural network where neural responses represent estimates of the latent state vectors and whose synaptic weights represent the spatial filters and transition parameters. We used K = 5 transition matrices for the first level (more matrices did not significantly improve performance—see Fig A in S1 Text). Perception in the DPC network corresponds to estimating the latent vectors by updating neural responses (through network dynamics) to minimize prediction errors via gradient descent (see Methods). Updating network parameters to further reduce prediction errors corresponds to learning (slow changes in synaptic weights through synaptic plasticity).

Fig 1d and 1e illustrate the inference process for both levels of the network. The network generates top-down and lateral predictions (green) using the current two-level state estimates

(blue). If the input sequence is predicted well by the top-down-modulated transition matrix V, the higher-level response  $\mathbf{r}^h$  remains stable due to small prediction errors (Fig 1d). When a non-smooth transition occurs in the input sequence, the resulting large prediction errors are sent to the higher level via feedforward connections (red arrows, Fig 1e, driving changes in  $\mathbf{r}^h$  to predict new dynamics for the lower level.

We trained the network on thousands of natural image sequences extracted from a video recorded by a person walking on a forest trail (frame size: 16 × 16 pixels, sequence length: 10 frames ( $\sim 0.35$  seconds)). The frames were spatially and temporally whitened to simulate retinal and lateral geniculate nucleus (LGN) processing [38, 41]. The image sequences reserved for testing did not overlap in space or time with the training sequences. Fig 2a illustrates the inference process on an example natural image sequence by the network. The first row displays the ground truth input  $I_t$  for 10 time steps: each frame was shown sequentially to the model. The next row shows the model's predictions  $U\bar{r}$ , for each time step t, where  $\bar{r}$ , was predicted by the previous state estimate  $\hat{\mathbf{r}}_{t-1}$ :  $\bar{\mathbf{r}}_t = \text{ReLU}(\mathbf{V}\hat{\mathbf{r}}_{t-1})$ . The prediction errors  $\mathbf{I}_t - \mathbf{U}\bar{\mathbf{r}}_t$  are shown in the third row. The prediction errors were the largest in the first two steps as the model inferred the spatial features and the transition dynamics from the initial inputs. The subsequent predictions were more accurate, resulting in minimized prediction errors. Finally, the last row shows the corrected estimates  $U\hat{\mathbf{r}}_t$  after  $\bar{\mathbf{r}}_t$  has been updated to  $\hat{\mathbf{r}}_t$  through prediction error minimization. Fig 2b shows the lower- (top) and higher-level (middle) neural responses to the natural video sequence in Fig 2a. The bottom panel of Fig 2b shows the topdown dynamics modulation generated by the higher level.

We examined the learned spatial receptive fields (RFs) of the model neurons at the first level and qualitatively compared them with the spatial RFs of simple cells in the primary visual cortex (V1). A subset of the spatial filters (columns of **U**) learned by the model from our natural videos dataset are shown in Fig 2c. These filters resemble oriented Gabor-like edge or bar detectors, similar to the localized, orientation-selective spatial RFs found in V1 simple cells [38, 42]. To measure the spatiotemporal receptive fields of the lower-level neurons, we ran a reverse correlation experiment [43, 44] with a continuous natural video clip ( $\approx$  47 minutes) extracted from the same forest trail natural video used for training. This video was not shown to the model during either training or testing (see Methods). Fig 2d shows the spatiotemporal receptive fields for four example lower-level model neurons, computed by weighting input frames from the seven previous time steps  $\mathbf{I}_{t-7}$ ,  $\mathbf{I}_{t-6}$ , ...,  $\mathbf{I}_{t-1}$  by the response  $\hat{\mathbf{r}}_t$  they caused at the current time step t (see Methods). The resulting average spatiotemporal receptive fields are shown as seven-image sequences labeled t-7, t-6, ..., t-1 (lasting  $\approx$  250 milliseconds in total). The first column labeled "Spatial" shows the spatial RFs of the example neurons.

To compute the space-time receptive fields (STRFs), we took the spatiotemporal X - Y - T receptive field cubes and collapsed either the X or Y dimension, depending on which axis had time-invariant responses. Fig 2e left panel shows the X/Y - T receptive fields of these example neurons. For comparison, Fig 2e right panel shows the STRFs of simple cells in the primary visual cortex (V1) of a cat (adapted from DeAngelis et al. [25]).

DeAngelis et al. [25] categorized the receptive fields of simple cells in V1 to be space-time separable (Fig 2e top row) and inseparable (Fig 2e bottom row). Space-time separable receptive fields maintain the spatial form of bright/dark-excitatory regions over time but switch their polarization: the space-time receptive field can thus be obtained by multiplying separate spatial and temporal receptive fields. Space-time inseparable receptive fields on the other hand exhibit bright/dark-excitatory regions that shift gradually over time, showing an orientation in the space-time domain. Neurons with space-time inseparable receptive fields are direction-selective, responding to motion in only one direction. As seen in Fig 2e left pane, the neurons in the lower level of our network learned V1-like separable and inseparable STRFs, based on the

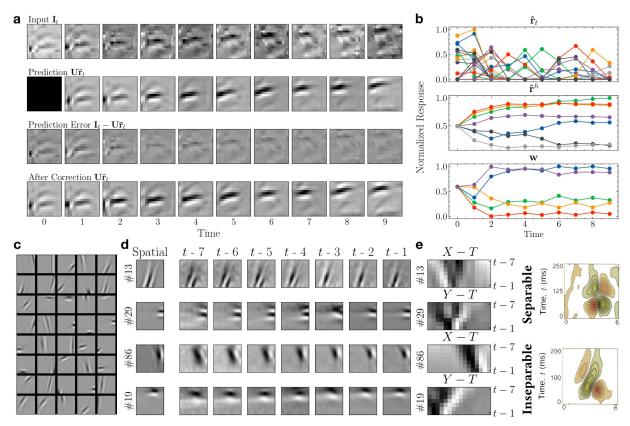


Fig 2. Predictive coding of natural videos and learned space-time receptive fields. (a) Inference on an example input image sequence of 10 frames. Top to bottom: Input sequence; model's prediction of the current input from the previous step (the first step prediction being zero); prediction error (predicted input subtracted from the actual input); model's final estimate of the current input after prediction error minimization. (b) The trained DPC network's response to the natural image sequence in (a). Each plotted line represents the responses of a model neuron over 10 time steps. Top: responses of the 20 most active lower-level neurons (some colors are repeated); middle: responses of seven randomly chosen higher-level neurons; bottom: predicted transition dynamics (each line is the modulation weight for a basis transition matrix at the lower level). (c) 40 example spatial receptive fields (RFs) learned from natural videos. Each square tile is a column of U reshaped to a  $16 \times 16$  image. (d) Space-Time RFs (STRFs) of four example lower-level neurons. First column: the spatial RFs of the example neurons. Next seven columns: the STRFs of the example neurons revealed by reverse correlation mapping. (e) Left panel: space-time plots of the example neurons in (d). Right panel: space-time plots of the RFs of two simple cells in the primary visual cortex of a cat (adapted from [25]).

principle of spatiotemporal prediction error minimization. To our knowledge, these results represent one of the first demonstrations of the emergence of both separable and inseparable STRFs in a recurrent network model by predictive coding of natural videos presented frame-by-frame. Previous demonstrations (e.g., [18, 20, 23]) have typically required chunks or all the frames of a video to be provided as a single input to a network, which is hard to justify biologically (see Discussion).

# Temporal hierarchy through prediction of dynamics

Next, we show that the two-level DPC network learned a hierarchical temporal representation of input videos. In our formulation of the model, a higher-level state vector predicts the *dynamics* of the lower-level states. This implies that the higher-level network neurons will have stable activation for input sequences with consistent dynamics (Fig 1d). When a change occurs in input dynamics, we expect the higher-level responses to switch to a different activation profile to minimize prediction errors (Fig 1e). We hypothesize that the different timescales of

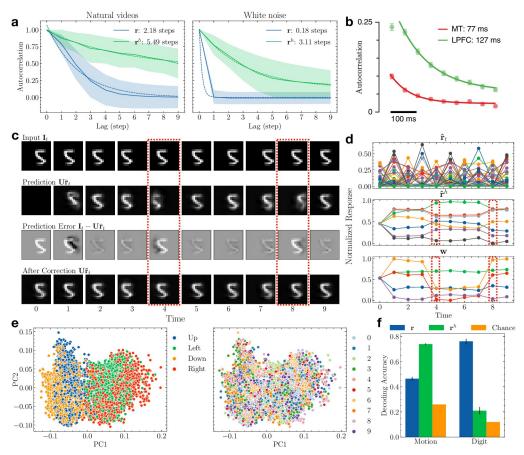
neural responses observed in the cortex [6–9] could be an emergent property of the cortex learning a similar hierarchical generative model.

We tested this hypothesis in our DPC network trained on natural videos. As seen in the inference example in Fig 2b, the lower-level responses change rapidly as the stimulus moves (top panel). The higher-level responses (middle panel) and the predicted transition dynamics (right panel) were more stable after the initial adaptation to the motion. Since the stimulus continued to follow roughly the same dynamics (leftward motion) after the first two steps, the transition matrix predicted by the higher-level neurons continued to be accurate for the rest of the steps, leading to small prediction errors and few changes in the responses. Note that we did not enforce a longer time constant or smoothness constraint for  $\mathbf{r}^h$  during inference—the longer timescale and more stable responses are entirely a result of the higher-level neurons learning to predict the lower-level dynamics under the proposed generative model.

To quantify this learned hierarchical temporal representation, we computed the autocorrelation of the lower- and higher-level responses to unseen natural videos and fitted an exponential decay function (see Methods). As Fig 3a shows, the autocorrelation of the higher-level responses  $\mathbf{r}^h$  is greater than that of the lower-level response  $\mathbf{r}$  and has a slower decay rate (exponential time constant for  $\mathbf{r}^h$ : 5.49 steps; for  $\mathbf{r}$ : 2.18 steps). To factor out the effect of the natural video statistics, we computed the same autocorrelation using Gaussian white noise sequences. We found that the hierarchy of timescales still exists ( $\mathbf{r}^h$ : 3.11 steps;  $\mathbf{r}$ : 0.18 steps). Fig 3b shows the autocorrelation of nonhuman primate neural responses from the medial-temporal (MT) area in the visual cortex (assumed to be lower in the processing stream) and lateral prefrontal cortex (LPFC) (assumed to be higher) for time periods preceding a motion task [6]. These neural responses show a difference in timescales qualitatively similar to the different response timescales in our hierarchical DPC network.

To further understand the model's ability to learn hierarchical temporal representations, we trained a DPC network on the Moving MNIST dataset [45]. Each image sequence in this dataset contains ten  $18 \times 18$  pixel frames showing a single example of a handwritten digit (chosen from the original MNIST dataset) moving in a particular direction. The digit's motion is limited to up, down, left, or right directions with a fixed speed. Fig 3c illustrates the trained network's inference process on an example image sequence. Similar to the responses to the natural video sequence, the lower-level responses displayed fast changes while the higher-level responses spanned a longer timescale and showed greater stability (Fig 3d). Note that at time t = 4 and t = 8, the input dynamics changed as the digit "bounced" against the boundaries and started to move in the opposite motion (Fig 3c red dashed box). The higher-level neurons' predictions resulted in large prediction errors at those times (Fig 3c third row). The prediction errors caused notable changes in the higher-level responses  $\mathbf{r}^h$  (Fig 3d red dashed boxes). For the rest of the steps,  $\mathbf{r}^h$  remained stable and generated accurate predictions of the stable dynamics.

Lastly, we confirmed that lower-level transition dynamics are indeed encoded in the higher-level responses. We performed principal component analysis (PCA) on the higher-level responses  $\mathbf{r}^h$  for the Moving MNIST sequences in the test set. Fig 3e visualizes these responses in the space of their first two principal components (PCs), colored by either the motion direction (left) or digit identities (right). The responses clearly formed clusters according to input motion direction but not digit identities. We then trained a support vector machine with radial basis function (RBF) kernel [46] to map  $\mathbf{r}$  and  $\mathbf{r}^h$  to motion directions and digit identities (Fig 3f). Using the higher-level responses, the classifier yielded 73.9% 10-fold cross-validated classification accuracy on the four motion directions (chance accuracy: 26.0%, computed as the number of majority labels in the test set divided by the total number of labels). Using the lower-level responses resulted in significantly less classification accuracy for motion direction



**Fig 3. Hierarchical temporal representation with different timescales. (a)** Autocorrelation of the lower- and higher-level responses in the trained network with natural videos. Shaded area denotes ±1 standard deviation. Dotted lines show fitted exponential decay functions. Left: response recorded during natural video stimuli; right: white noise stimuli. **(b)** Autocorrelation of the neural responses recorded from MT and LPFC of monkeys. Adapted from Murray et al. [6] **(c)** Inference for an example Moving MNIST sequence in a trained network. The red dashed boxes mark the time steps when the dynamics of the input changed. **(d)** The network's responses to the input Moving MNIST sequence in (c). Note the changes in the higher-level responses after the input dynamics changed (red dashed boxes); this gradient-based change helps to minimize prediction errors. **(e)** Higher-level responses to the Moving MNIST sequences visualized in the 2D space of the first two principal components. Left: responses colored according to motion direction; right: responses colored according to digit identities. **(f)** Comparison of decoding performance for motion direction versus digit identity using lower- and higher-level neural responses. Error bars: ±1 standard deviation from 10-fold cross validation. Orange: chance accuracies.

(46.5%,  $p \ll 0.001$ , t-test). In contrast, decoding accuracy for digit identity was significantly higher using the lower-level responses (76.1%) compared to using the higher-level responses (20.9%,  $p \ll 0.001$ , t-test). These results show that due to the structure of its generative model, the DPC network learned to disentangle to a significant extent the motion information in an input video from image content (here, digit identity), yielding a factored representation of input image sequences.

#### Predictive and postdictive effects in visual motion processing

The ability of the DPC model to encode entire sequences at the higher level (cf. the "timeline" model of perception [29]) leads to new normative and computational interpretations of visual motion phenomena such as the flash-lag illusion [26, 30, 31], explaining both predictive and

postdictive effects [27, 29]. The flash-lag illusion refers to the phenomenon that a flashed, intermittent object is perceived to be "lagged" behind the percept of a continuously moving object even though the physical locations of the two objects are aligned or the same [30, 31]. Though this illusion is commonly attributed to the predictive nature of the perceptual system [30], Eagleman and Sejnowski [26] proposed a postdictive mechanism based on psychophysical results that the motion of the moving object *after* the flash can change the percept of events at the time of the flash. The potential interplay between prediction and postdiction in shaping perception was also studied by Hogendoorn et al. [27, 29]. The authors designed an interference paradigm with different reaction-speed tasks and showed that when the trajectory of the object unexpectedly reverses, predictive effects (extrapolation) are observed at short latencies but postdictive effects (interpolation) are observed at longer latencies (Fig 4i and 4i).

We propose that prediction error minimization with a hierarchical temporal representation, as in the DPC model, provides a natural explanation for these predictive and postdictive effects. In a DPC network, the higher-level state  $\mathbf{r}^h$  predicts entire sequences of lower-level states following the same dynamics (Fig 3). When the dynamics of observations change (e.g., motion reversal), the higher-level state is updated to minimize prediction errors, resulting in a revised state that represents the motion-reversed sequence spanning both past and future inputs. This process corresponds to postdiction in visual processing [28]. For the flash-lag experiment, we predict that the higher-level neurons of a trained DPC network will form a static sequence percept when presented with a flashed object and a directional sequence percept for a moving object, causing perceived lags between the two objects as observed in the flash-lag illusion [30].

We first test these predictions of the DPC model on the experimental conditions used by [26]. In their experiment, the stimuli consisted of a flashed disk and a ring moving in a circle. Before the flash, the ring could have an initial trajectory (Fig 4a top) or no initial trajectory (Fig 4a bottom). After the flash, the ring could continue moving on its initial trajectory ("continuous"), stop moving ("stopped"), or move on the reversed trajectory ("reversed"). A flash appeared in a seven-degree range that extended above and below the ring on its trajectory. The participants then indicated whether a flashed white disk occurred above or below the center of the moving ring. Positive displacements denoted lags along the initial trajectory of the ring, while negative displacements denoted the reversed direction.

To simulate these testing conditions, we used the Moving MNIST test set and extracted 134 test sequences with consistent leftward or rightward motion. For each of these 134 sequences, we simulated the two test conditions used by [26] (with or without initial trajectory): the higher-level state  $\hat{\mathbf{r}}^h$  was either inferred from the first three steps (t = 0, 1, 2) of the input sequence, or initialized to the zero vector (Fig 4b left). For each of these two test conditions, we simulated the four test cases used in [26] regarding the motion of the moving object at the time of the flash (Fig 4b right). Note that flashed stimuli correspond to the "no initial trajectory, terminate" condition. We computed the location of a digit as the center of mass of pixel values in the 2D image; the perceived location at time t was defined similarly based on the predicted image at time t. As Fig 4e shows, the perceived location of a flashed object at t = 3strongly overlapped with the physical flashed location at t = 2, showing that the prediction errors drove the higher-level state estimates to predict no change in object location for the flashed object. Fig 4f shows the perceived displacement between the moving object (with initial trajectories) and the flashed object, computed as the difference in perceived locations at t = 3between the moving object and the flashed object. Positive displacements followed the original trajectory direction and negative displacements followed the reversed direction. The perceived displacements in the model were significantly different in the three test conditions (Fig 4f left panel,  $p \ll 0.001$ , one-way ANOVA test) and were similar to the psychophysical results

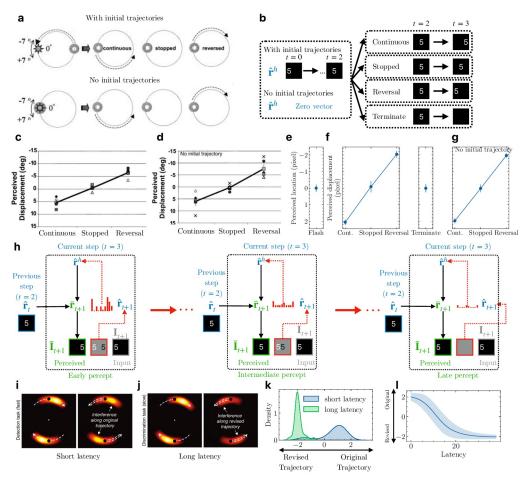


Fig 4. Flash-lag illusion and object representations in apparent motion. (a) The flash-lag test conditions used by [26]. The moving ring could have an initial trajectory (top) or no trajectory (bottom). At the time of the flash (bright disk), the ring could move along the initial trajectory, stop, or reverse its trajectory. Adapted from [26]. (b) Two test conditions (left) regarding initial trajectories of the moving object (a digit) in the flash-lag experiment with the model, and four test conditions (right) for the moving object. The flashed object was shown at time t and turned off at time t + 1 (same as the "Terminate" condition). (c & d) Psychophysical estimates for human subjects reported by [26] when the moving object had initial trajectories (c) or no initial trajectory (d). (e) Perceived location of the flashed object in the DPC model at time t + 1. The error bar indicates  $\pm 1$  standard deviation (measured across presentations of different digits). (f) Perceived displacement between the moving object (with initial trajectories) and the flashed object in the DPC model for the four test conditions. (g) Same as (f) but with no initial trajectory for the moving object. (h) Illustration of the predictionerror-driven dynamics of the perception of the moving object in the model when the trajectory reversed at time t + 1. Red ellipsis between panels denotes the prediction error minimization process. (i) Interference pattern during human apparent motion perception with continuous motion (left) and reversed motion (right) at short latency (fast detection task). Brighter color denotes more interference. Dashed arrows represent object motion direction. Adapted from [29]. (j) Same as (i) but at long latency (slow discrimination task) [29]. (k) Perceived location of the moving object in the DPC  $model\ at\ time\ t+1\ probed\ at\ short\ versus\ long\ latency\ during\ prediction\ error\ minimization.\ Positive\ values\ denote$ distance along the original trajectory. Negative values denote distance along the reversed trajectory. Short and long latency correspond to "Early percept" and "Late percept" respectively in part (h). (I) Perceived location of the digit at all latencies during the prediction error minimization process in part (h).

reported by Eagleman & Sejnowski (Fig 4c). Fig 4g confirms that the initial trajectories of the moving object had no effects on the model's flash-lag illusion, consistent with the reported results (Fig 4d) [26]. These results validate the explanation provided by the DPC model on the flash-lag effect: for a hierarchical generative model with representations of sequences, a flashed or stopped/terminated moving object leads to inference of a static object sequence (Fig 4e),

while continuous or reversed motion leads to inference of a moving object sequence, resulting in the perceived lags along the corresponding directions (Fig 4f).

One aspect of motion perception the previous results do not illustrate is the interplay between postdiction and prediction. Hogendoorn et al. investigated this effect in an experiment on apparent motion perception [27]. Participants were instructed to report the detection of a visual cue (short latency task) or differentiate between two visual cues (long latency task) during apparent motion. These visual cues could either be along the apparent motion trajectory or the reversed trajectory. The authors found that upon reversing the apparent motion trajectory, predictive effects dominated perception at short latency (detection task, Fig 4i), with the most interference (measured in terms of the participants' reaction times) along the original motion trajectory. At longer latency (differentiation task, Fig 4j), most interference was along the reversed trajectories, indicating that postdictive effects dominated perception.

We hypothesize that the prediction error minimization process of DPC could explain this interplay between prediction and postdiction, as illustrated by Fig 4h which depicts the gradient-descent-based optimization process of Fig 1e (and Eq 18). Early percepts of the model are dominated by the spatiotemporal prediction using the optimal estimates from the previous step (Fig 4h left). When a motion reversal occurs, feedforward prediction errors gradually correct the second-level state (Fig 4h middle) until convergence (Fig 4h right). Therefore, late percepts in the model correspond to error-corrected spatiotemporal predictions. Note that due to the discrete temporal nature of the DPC model (unit time steps), this process is considered to happen "at" one particular time step (e.g., early versus late percept "at" t = 3 in Fig 4h).

To test this hypothesis, we used the same trained DPC network and probed its percept of the moving object at the time of reversal under the "with initial trajectory, reversal" condition (Fig 4b). At short latency (10% of steps into prediction error correction, Fig 4h early percept), the perceived locations for the moving object in most test sequences were along the original trajectory, as denoted by positive displacements compared to the final step before reversal (t = 2) (Fig 4k blue)). At longer latency (90%, Fig 4h late percept), the moving object's perceived locations were flipped and along the reversed trajectory (negative displacements; Fig 4k green,  $p \ll 0.001$ , t-test). This is consistent with psychophysical findings [27, 29] that when the motion of the object unexpectedly reversed, prediction effects were observed at short latency ( $\approx$  350 ms, Fig 4i right panel, bright color denotes locations of interference due to prediction) while postdictive effects were observed at longer latency ( $\approx$  620 ms, Fig 4j right panel, bright color denotes locations of interference due to postdiction). Fig 4l plots the moving object's perceived location in our model throughout the error correction process: the perceived location varies smoothly from being along the original direction initially to along the reversed direction at greater latencies. These results make a testable prediction: if probed at an intermediate level of latency (between 350 ms and 620 ms), the maximal interference should overlap with the object's location at the time of reversal (i.e., at the black dots in Fig 4i and 4j), as suggested by Fig 41.

## Cue-triggered recall and episodic memory

A number of experiments in rodents have shown that the primary visual cortex (V1) encodes predictive representations of upcoming stimuli [1–4, 47]. In one of the first such studies, Xu et al. [1] demonstrated that after exposing rats repeatedly to the same moving dot visual sequence (Fig 5a), displaying only the starting dot stimulus triggered sequential firing in V1 neurons in the same order as when displaying the complete sequence (Fig 5a). Similar effects have been reported in monkey [35] and human [32–34, 48] visual cortical areas as well.

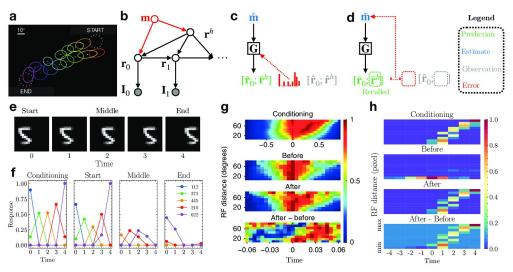


Fig 5. Cue-triggered activity recall in the DPC model. (a) The experimental setup of Xu et al. (adapted from [1]). A bright dot stimulus moved from START to END repeatedly during conditioning. Activities of neurons whose receptive fields (colored ellipses) were along the dot's trajectory were recorded. (b) Generative model combining an associative memory and DPC. The red part denotes the augmented memory component that binds the initial content vector  $\mathbf{r}_0$ and the dynamics vector  $\mathbf{r}^h$  to encode an episodic memory. (c) Depiction of the memory encoding process. The presynaptic memory activity and postsynaptic prediction error jointly shape the memory weights G. (d) Depiction of the recall process. Prediction error on the partial observation  $\hat{\mathbf{r}}_0$  drives the convergence of the memory estimates  $\tilde{\mathbf{m}}$ and recalls the higher-level dynamics vector  $\bar{\mathbf{r}}^h$  as a top-down prediction. The red dotted box depicts the prediction error between the missing observations for  $\mathbf{r}^h$  and the prediction  $\bar{\mathbf{r}}^h$ ; this error is ignored during recall, implementing a form of robust predictive coding [49]. (e) The image sequence used to simulate conditioning and testing for our memory-augmented DPC network. (f) Responses of the lower-level neurons of the network. Colored lines represent the five most active lower-level neurons at each step. Left to right: neural responses during conditioning, testing the network with a single start frame, middle frame, and end frame. (g, h) Normalized pairwise cross correlation of (g) primary visual cortex neurons (adapted from [1]) and (h) the lower-level model neurons. Top: during conditioning; middle two: testing with the starting stimulus, before and after conditioning; bottom: the differences between cross correlations, "After" minus "Before" conditioning.

The generative model of DPC provides a highly efficient computational basis for episodic memories and sequence prediction. DPC assumes sequences are generated by a factorized representation: a single (lower-level) representation of the content ("what") provided at the first step and a single (higher-level) representation of dynamics (motion or "where"). These two representations are inferred during sequence perception as the explanations (or causes) of a given input sequence.

It is known that factored information from the visual cortex makes its way, via the medial and lateral entorhinal cortices, to the hippocampus [50]. The hippocampus has been implicated both in the formation of episodic memories [51–53] and in mediating activity recall in the neocortex [54–57] through its outputs to the entorhinal cortex, which in turn conveys this information to downstream areas via feedback connections. Because the DPC model encodes an entire sequence in terms of a single dynamics vector  $\mathbf{r}^h$  (along with the content  $\mathbf{r}_0$  at the first step), it suggests a simple mechanism for storing sequential experiences as episodic memories, namely, storing the vector  $\mathbf{r}^h$  (along with  $\mathbf{r}_0$ ) in an associative memory, mimicking the role of the hippocampus.

To test this hypothesis, we augmented the DPC model with an associative memory that uses a vector  $\mathbf{m}$  to bind the content vector  $\mathbf{r}_0$  with the dynamics of the sequence  $\mathbf{r}^h$ , thereby encoding an episodic memory: the new generative model is shown in Fig 5b. Given the initial cue  $\mathbf{r}_0$  (inferred from the first image frame in the sequence), the associative memory

(emulating the hippocampus) recalls the episodic sequence dynamics  $\mathbf{r}^h$  which modulates the transition dynamics in the DPC network (representing the visual cortex) to complete the sequence recall. Specifically, we added to the trained DPC network another higher level of predictive coding to implement associative memory [10, 58]: the memory vector estimate  $\hat{\mathbf{m}}$  predicts both the content vector  $\hat{\mathbf{r}}_0$  and motion dynamics vector  $\hat{\mathbf{r}}^h$  and uses the prediction error to correct itself (Fig 5c). Upon convergence of  $\hat{\mathbf{m}}$ , the associative memory network stores this vector by updating its weights  $\mathbf{G}$  using Hebbian plasticity based on presynaptic activity  $\hat{\mathbf{m}}$  and postsynaptic prediction error [10] (Fig 5c, see Methods). During recall, the memory estimate  $\hat{\mathbf{m}}$  is driven by the  $\hat{\mathbf{r}}_0$  inferred from the cue and the prediction error (Fig 5d, dashed boxes denote the missing input). The dynamics vector  $\mathbf{r}^h$  is then recalled as the top-down prediction after  $\hat{\mathbf{m}}$  has converged (Fig 5d, green dashed box). Note that during conditioning, no learning occurs in the DPC network—only the weights of the memory network  $\mathbf{G}$  are optimized to store the episodic memory  $\hat{\mathbf{m}}$ .

We simulated the experiment of Xu et al. [1] using a moving MNIST sequence from the test set shown in Fig 5e. After conditioning (5 repetitions of the sequence), the network was tested with the starting frame only, the middle frame only, and the end frame only. The lower-level responses  $\hat{\mathbf{r}}_0$  of the DPC network were used to recall the dynamics component  $\bar{\mathbf{r}}^h$  from the memory. The recalled dynamics were then used to predict a sequence of lower-level responses in the DPC network. We found that the lower-level model neurons exhibited cue-triggered activity recall given only the start frame of the sequence (Fig 5f Start). Cueing the network with the middle frame triggered weak recall, consistent with findings by Xu et al. (see Fig 3c in Ref [1]). The end frame did not trigger recall [1]. We found that the sequence recall is cue-specific —when trained with sequences that have distinct digits and dynamics, the DPC network successfully recalled the correct sequence when cued with different starting digits (Fig B in S1 Text).

Lastly, following the analysis done by Xu et al. [1], we plotted the pairwise cross-correlation of the lower-level model neurons as a function of their spatial RF distances when tested with the starting frame of the sequence (see Methods). As Fig 5h shows, the peaks of the correlation showed a clear rightward slant after conditioning, consistent with the experimental results (Fig 5g). This indicates a strong sequential firing order in the lower-level model neurons elicited by the starting cue, where neurons farther apart have longer lags in response cross-correlations, a phenomenon that was nonexistent before conditioning (Fig 5g). These simulation results support our hypothesis that cue-triggered recall could be the result of the hippocampus, acting as an associative memory, binding factorized sequence representations of content and dynamics from the neocortex and recalling the corresponding dynamics component given the content cue.

# Estimating higher-order transition dynamics with a three-level model

Our results thus far involved a two-level DPC model whose second-level states predicted the first-level state transitions. Since the second-level state is assumed to characterize the *entire* sequence (Fig 1), it cannot predict higher-order transitions such as digits bouncing at the boundaries in the Moving MNIST dataset, which requires different second-level state representations (Fig 3). Here we show that adding a third level allows the DPC model to learn and infer the *transition dynamics of second-level states*, thus capturing a temporally more abstract representation of sequences at the third level.

Fig 6a shows the generative model for the three-level DPC model. Just as the second-level states modulate the transition function of the first-level states in the two-level model (Fig 1), the third-level states modulate the transitions of the second-level states. During inference (Fig

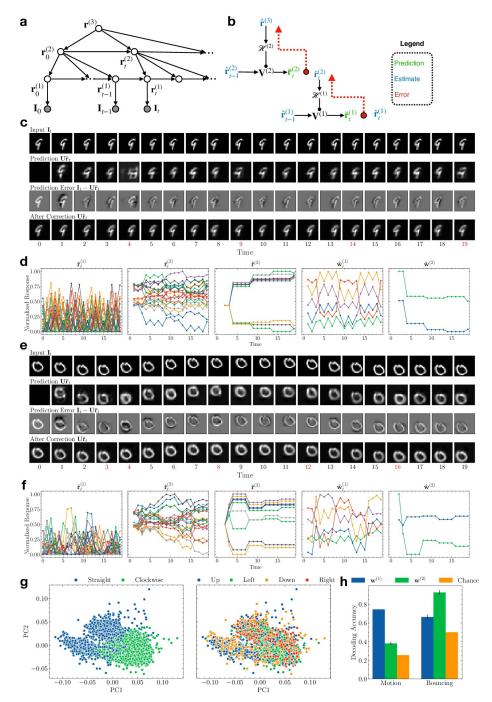


Fig 6. Three-level DPC model learns progressively more abstract temporal representations. (a) Generative model for three-level DPC. (b) Schematic depiction of an inference process. Observation nodes are omitted for clarity. (c) Inference for an example Moving MNIST sequence with "straight bouncing" dynamics. Red time steps mark the moments when the first-level prediction error exceeded the threshold, causing the network to transition to a new second-level state (see Methods). For these time steps, the predictions (second row) are by the second-level neurons, while the rest are by the first-level neurons as in Fig 3. (d) The network's responses to the Moving MNIST sequence in (c). Left to right: first-level responses, second-level responses, third-level responses, first-level modulation weights, second-level modulation weights. (e) Same as (d) but with "clockwise bouncing" dynamics. (f) Same as (d) but for the sequence in (e). (g) Third-level responses to the Moving MNIST sequences visualized in the 2D space of the first two principal components. Left: responses colored according to bouncing type; right: responses colored according to motion direction. (h) Comparison of decoding performance for bouncing type versus motion direction using the modulation weights generated by the second and third level. Error bars: ±1 standard deviation from 10-fold cross validation. Orange: chance accuracies.

6b), when the first-level prediction error is larger than a threshold (see Methods), the second-level state transitions to the next state, following the transition function predicted by the current third-level state. The second-level prediction error is conveyed to the third level to correct its state estimate, in the same way as the first-level prediction error corrects the second-level state.

The Moving MNIST dataset we used for Fig 3 exhibited only one type of transition dynamics of the second-level states, namely, transitioning from moving left to moving right, or moving up to moving down, and vice versa (henceforth referred to as "straight bouncing" dynamics (Fig 6c). To demonstrate that the third-level states can learn different second-level dynamics, we added to the dataset digit sequences with "clockwise bouncing" dynamics to the dataset (for example, a digit moving to the left and hitting the boundary will move upward instead of rightward, and so on (Fig 6e). This makes the second-level state transition function ambiguous until the first bouncing event. If the third-level representations learned by DPC capture the second-level transition dynamics, we expect the first large prediction error at the second level (occurring at a boundary) to update the third-level state estimate to represent either straight bouncing dynamics or clockwise bouncing dynamics. Thereafter, the third-level state estimate should remain stable as long as the bouncing type remains the same.

In the following, we use the superscript (i) to denote the level i. We trained a three-level neural network on the augmented Moving MNIST dataset (with the two types of bounding dynamics discussed above). The network uses two second-level transition matrices  $\{\mathbf{V}_1^{(2)}, \mathbf{V}_2^{(2)}\}$ and a top-down network  $\mathcal{H}^{(2)}$  (from the third to the second level), in addition to all the parameters in the two-level model. The first and second-level transition matrices were pretrained (see Methods). Fig 6c and 6d show an inference example of the three-layer network on an input sequence with straight bouncing dynamics. Red time steps denote the moments when the first-level prediction errors were larger than the set threshold, causing the second-level neurons to change their activities to transition to the next state (this can be seen as a neural implementation of terminal states in a hierarchical HMM [59] (see Methods)). As seen in the second row in Fig 6c, at the first bouncing event (t = 4), the second-level prediction was not accurate; the third-level neural responses were updated to minimize this prediction error (see Fig 6d). For the rest of the sequence, the predictions are accurate at the bouncing events (t = 9, 14, 19) and the third-level neural responses remained stable. The panels in Fig 6d show an increase in response stability and timescale from the first to the third-level neural responses (first three panels), as well as in the modulation weights that define the first and second-level transition dynamics (last two panels). Fig 3e and 3f show a different example with clockwise bouncing dynamics. Similar to the example above, the third-level responses showed notable changes at times t = 3 and 4 but remained stable for the rest of the sequence. Comparing the second-level modulation weights in Fig 6d and 6f, it is clear that the third-level DPC neurons estimated different bouncing types and generated opposite modulation strengths for the two types of sequences.

We performed PCA on the third-level responses  $\mathbf{r}^{(3)}$  obtained at the end of the Moving MNIST sequences (t=19) in the test set. Fig 6g visualizes these responses in the space of their first two principal components (PCs), colored either by bouncing dynamics type (left) or moving direction (right). The third-level responses form clusters according to the bouncing dynamics type but not motion direction. We then used SVMs to decode the bouncing dynamics type or motion direction from  $\mathbf{w}^{(1)}$  and  $\mathbf{w}^{(2)}$ , the weights predicted by  $\mathbf{r}^{(2)}$  and  $\mathbf{r}^{(3)}$  respectively. As shown in Fig 6h, using  $\mathbf{w}^{(2)}$ , the classifier yielded 92.7% 10-fold cross-validated classification accuracy on the two bouncing types (chance accuracy: 50.0%). Using  $\mathbf{w}^{(1)}$  resulted in significantly less classification accuracy for bouncing type (66.4%,  $p \ll 0.001$ , t-

test). In contrast, decoding accuracy for the four motion directions (chance accurary: 25.4%) was significantly higher using  $\mathbf{w}^{(1)}$  (74.5%) compared to using  $\mathbf{w}^{(2)}$  (38.0%,  $p \ll 0.001$ , t-test). These results show that the three-level DPC model succeeded in learning a temporal hierarchy, with the third-level states encoding the longest timescale feature, *i.e.* the type of bouncing dynamics, by modulating the transition function of the second-level states, which in turn encoded intermediate timescale features (motion direction).

#### **Discussion**

Our results show that dynamic predictive coding (DPC) can learn hierarchical temporal representations of sequences through top-down modulation of lower-level dynamics. Specifically, we showed that by minimizing prediction errors on image sequences, a two-level DPC neural network develops V1-like separable and inseparable space-time receptive fields at the lower level [25], and representations encoding sequences at a longer timescale at the higher level [6, 8]. The trained DPC network provides a normative explanation for the flash-lag effect [26] and accounts for both prediction and postdiction in visual motion processing [27, 28, 60]. The temporal abstraction of sequences in a DPC network suggests a new mechanism for storing and retrieving episodic memories by linking the DPC network to an associative memory, emulating the interaction between the neocortex and the hippocampus. We show that such a memory-augmented DPC model explains cue-triggered activity recall in the visual cortex [1]. Finally, we show that the top level of a three-level DPC network captures the higher-order temporal statistics encoding the transition dynamics of the second-level states, which in turn capture the temporal statistics of the first-level states. Taken together, the hierarchical temporal representations learned by DPC, ranging from the lowest-level space-time representations similar to those observed in visual cortical simple cells (Fig 2), through the intermediate-level representations of steady motion (Fig 3), to the highest-level representation of how such motion changes over a longer timescale (Fig 6), emulate the spatiotemporal representations observed in visual cortical hierarchies, particularly along the dorsal visual pathway [61].

The key to the DPC model's ability to capture lower-level dynamics with relatively stable higher-level response vectors  $\mathbf{r}^h$  is the top-down modulation of transition dynamics of entire lower-level state sequences, using the weights w generated by the higher level. There has been increasing interest in neuroscience in the role of modulatory inputs (e.g., encoding top-down contextual information) in shaping the dynamics of recurrent neural networks in the brain [62-64]. The DPC model ascribes an important role to these modulatory inputs in enabling cortical circuits to learn temporal hierarchies. The neural implementation used in this paper can be seen as top-down feedback (w) targeting the distal apical dendrites of lower-level pyramidal neurons, thereby changing their gain (Fig 1c). Such a mechanism, which has been shown to be possible experimentally [16, 17, 39, 65], can also modulate perceptual detection [66]. Note that although we chose to model the top-down influence as multiplicative gain modulation, it would be theoretically equivalent to model it as an additive component or concatenate it as an extra input for prediction (e.g., predicting first-level transitions as  $\bar{\mathbf{r}}_t = \mathbf{f}(\mathbf{r}_{t-1}, \mathbf{r}^h)$ , where  $\mathbf{f}$  is a multi-layer perception). However, such an implementation may be less efficient (in terms of the number of parameters required to reach the same level of performance) under certain conditions, compared to a hypernetwork-based implementation [67] such as our implementation based on multiplicative gain modulation.

Some of the first models of spatiotemporal predictive coding focused on signal processing in the retina and LGN [41, 68]. Other models for sequence processing, such as sparse coding [20, 38] and independent component analysis [18], have been shown to produce oriented space-time receptive fields from natural image sequences, but these models require the entire

image sequence to be presented as a single vector input, which is hard to justify biologically; they also do not explicitly model the temporal dynamics between images and therefore, cannot make predictions into the future given a single input. A previous spatiotemporal predictive coding model based on Kalman filtering [40] did incorporate state transitions but the model was not hierarchical and was not shown to generate cortical space-time receptive fields. Our model bears some similarities to slow feature analysis which extracts slowly varying features from sequences of stimuli but it does not learn the transition dynamics between time steps [19, 21, 22]. DPC on the other hand learns a generative model that generates entire sequences, with the assumption that the transition dynamics do not change within a sequence (a "slow" feature). Object identity remains in the lower-level representations of DPC (Fig 3f). From a learning perspective, Luczak et al. [69] propose that single neurons predicting their future activity at a fixed delay could also serve as an effective learning mechanism.

Recent advances in deep learning [70] have spurred several efforts to learn spatiotemporal hierarchies from sensory data. Lotter et al. developed a deep learning model called "PredNet" for learning a hierarchical predictive coding-inspired model for natural videos [71, 72]. After training, the model was shown to produce a wide range of visual cortical properties and motion illusions. However, in PredNet, higher-level neurons predict lower-level prediction errors rather than neural activities or dynamics, making it unclear what the underlying generative model is. It is also unclear if PredNet learns a temporal response hierarchy as found in the cortex. A different model, proposed by Singer et al. [23] and later extended to hierarchies [24], is trained by making higher layers predict lower layer activities: after training, model neurons in different layers displayed different levels of tuning properties and direction selectivity similar to neurons in the dorsal visual pathway. However, similar to the sparse coding and ICA models discussed above for spatiotemporal sequences, the Singer et al. model also requires a sequence of images to be presented as a single input to the network, and the hierarchy of timescales is hard-coded (higher-level neurons predict future lower-level neural activities by receiving a fixed-length chunk of neural activities as input). The above models also do not provide explanations for postdiction or episodic memory and recall.

Many experimental studies have shown an increase in temporal representation stability and response timescales as one goes from lower-order to higher-order areas in the visual and other parts of the cortex [6–9, 73, 74]. Most computational models have studied this phenomenon through mechanistic rate-based models with parameters based on connectivity data [75, 76] or spiking network models [77]. Kiebel et al. [78] proposed a model where second-level states generate a single parameter for the first-level Lorenz attractor as the slower "sensory cause" parameter. DPC generalizes this model by assuming higher-level states fully determine the lower-level transition function by predicting the transition dynamics of lower-level states. Under this formulation, temporal hierarchies emerge naturally as a consequence of the neocortex learning from temporally structured data (e.g., stable dynamics in short time windows). This view is consistent with findings that response timescales are functionally dynamic and could expand for cognitive tasks such as working memory [79].

Previous normative models of postdiction in visual processing often relate the effect to the concept of Bayesian smoothing (or backward message passing) [26, 80]. We have shown that a trained two-level DPC network with higher-level sequence representations also exhibits postdictive effects without the need for smoothing. In the event of a temporal irregularity (e.g., an unexpected motion reversal), the higher-level state in the DPC network is updated to reflect a new revised input sequence, naturally implementing postdiction through online hierarchical Bayesian filtering (Figs 1 and 3). Our flash-lag simulation results are consistent with the Bayesian filtering model from Khoei et al. [81] showing that the flash-lag effect can be produced through an internal model that explicitly represents object velocity. The higher-level sequence

representation in the DPC model supports an implicit (and more generalized) representation of velocity and reproduces the same internal dynamics of the "speed" estimate at motion reversal (compare Fig 4h with Fig 6 in [81]). It is worth noting that the trained DPC network learned to predict no motion (static sequence) for the flashed object even though it was never trained on static object sequences and did not assume a prior of zero speed [81]. This emergent property was also seen in PredNet, which learned to predict relatively little motion for a flashed bar stimulus [72].

The higher-level sequence representations of DPC, when combined with an associative memory, support the formation of episodic memories and cue-triggered activity recall [1, 3, 32–35]. The associative memory in our model forms an episodic memory by binding the inferred content representation and dynamics representation from the DPC network during conditioning. When an initial portion of the sequence is presented during testing, the stored episodic memory is retrieved, generating the dynamics component which modulates the lower-level network to enable full recall of the sequence. Though previously considered to only require V1 plasticity [1, 3], sequence learning is severely impaired in mice with hippocampal damage [82]. Coordinated activity between V1 and the hippocampus has also been found in human V1 during recall [48]. These experimental results support the involvement of the hippocampus in sequence learning, consistent with our model. Overall, the memory-augmented DPC model offers a highly efficient computational basis for forming and recalling episodic memories [57, 83], where a single representation of content and transition dynamics from all sensory areas of the neocortex can be bound together as a memory and later retrieved upon receiving partial input.

In our three-level DPC model, the second-level state (under the influence of the current third-level state) predicts the next second-level state only when the first-level prediction errors are larger than an estimated threshold. This can be seen as a neural (and continuous-valued) implementation of "terminal states" in hierarchical hidden Markov models (HHMMs) [59]. In an HHMM, when a terminal state is reached at a lower level, the corresponding sub-HMM is deemed to be completed and the higher-level state then transitions to the next higher-level state (which activates the next sub-HMM at the lower level). In our hierarchical DPC model, small first-level prediction errors are resolved locally between the first and second level, indicating a continuing sub-sequence. When the error exceeds the threshold, the sub-sequence ends and the second-level transitions are activated. Any second-level prediction errors are resolved between the second and third level through third-level state inference. We used posthoc estimation of the error threshold after training but future work could attempt to estimate the threshold online in terms of the inverse variance or "precision" of prediction errors [84]). Additionally, second-level transitions could also correlate with the spatial information in the videos (e.g. bouncing only happens when the digit is near the boundary). Models whose second-level states depend on both the previous first- and previous second-level states could learn this type of transition [85].

The DPC model can be extended to action-conditioned prediction and hierarchical planning (see, e.g., [86] for initial steps in this direction). There is a growing body of evidence that neural activity in the sensory cortex is predictive of the sensory consequences of an animal's own actions [2, 4, 13, 47, 87]. These results can be understood in the context of a DPC model in which the transition function at each level is a function of both a state and an action at that level, thereby allowing the hierarchical network to predict the consequences of actions at multiple levels of abstraction [86]. Such a model allows probabilistic inference to be used not only for perception but also for hierarchical planning, where actions are selected to minimize the sensory prediction errors with respect to preferred goal states. Such a model is consistent with theories of active inference [11] and planning by inference [88–92], and opens the door to

understanding the neural basis of navigation and planning [9, 93, 94] as an emergent property of prediction error minimization.

#### **Methods**

## Hierarchical generative model

We assume that the observation  $\mathbf{I}_t \in \mathbb{R}^M$  at time t is generated by a lower-level latent variable  $\mathbf{r}_t \in \mathbb{R}^N$ . The latent variable  $\mathbf{r}_t$  is generated by the previous step latent variable  $\mathbf{r}_{t-1}$  and the higher-level latent variable,  $\mathbf{r}^h \in \mathbb{R}^{N_h}$ . Together, the generative model factorizes as follows:

$$p(\mathbf{I}_{0:T-1}, \mathbf{r}_{0:T-1}, \mathbf{r}^h) = p(\mathbf{r}^h)p(\mathbf{r}_0) \prod_{t=0}^{T-1} p(\mathbf{I}_t \mid \mathbf{r}_t) \prod_{t=1}^{T-1} p(\mathbf{r}_t \mid \mathbf{r}_{t-1}, \mathbf{r}^h).$$
(3)

Each component of the factorization is parameterized as follows:

$$\mathbf{r}^h \sim \mathcal{N}(0,1) \tag{4}$$

$$\mathbf{r}_{t} \mid (\mathbf{r}_{t-1}, \mathbf{r}^{h}) \sim \mathcal{N}(\mathbf{f}(\mathbf{r}^{h}, \mathbf{r}_{t-1}), \sigma_{r}^{2} I)$$
(5)

$$\mathbf{I}_t \mid \mathbf{r}_t \sim \mathcal{N}(\mathbf{U}\mathbf{r}_t, \sigma^2 I),$$
 (6)

where  $\mathcal{N}$  denotes the normal distribution and I denotes the identity matrix. The mean  $\mathbf{r}_t^{\mu} = \mathbf{f}(\mathbf{r}^h, \mathbf{r}_{t-1})$  is given by:

$$\mathbf{w} = \mathcal{H}_{\theta}(\mathbf{r}^h) \tag{7}$$

$$\mathbf{V} = \sum_{k=1}^{K} w_k \mathbf{V}_k \tag{8}$$

$$\mathbf{r}_{\star}^{\mu} = \text{ReLU}(\mathbf{V}\mathbf{r}_{t-1}). \tag{9}$$

Here,  $\mathcal{H}_{\theta}$  is a function (neural network) parameterized by  $\theta$ .

To sum up, the trainable parameters of the model include spatial filters  $\mathbf{U}$ , K transition matrices  $\mathbf{V}_1, \ldots, \mathbf{V}_K$ , and the neural network parameters  $\theta$ . The latent variables are  $\mathbf{r}_{0:T-1}$  and  $\mathbf{r}^h$ . See "Summary of the DPC generative model" in S1 Text for a more detailed description of the model architecture.

## **Prediction error minimization**

Here, we derive the loss function used for inference and learning under the assumed generative model. We focus on finding the *maximum a posteriori* (MAP) estimates of the latent variables using a Bayesian filtering approach. At time t, the posterior of  $\mathbf{r}_t$  conditioned on the input observations up to time t,  $\mathbf{I}_{0:t}$ , and the higher-level variable  $\mathbf{r}^h$  can be written as follows using Bayes' theorem:

$$p(\mathbf{r}_{t} \mid \mathbf{I}_{0:t}, \mathbf{r}^{h}) \propto p(\mathbf{I}_{t} \mid \mathbf{r}_{t})p(\mathbf{r}_{t} \mid \mathbf{I}_{0:t-1}, \mathbf{r}^{h}), \tag{10}$$

where the first term on the right-hand side is the likelihood function defined by Eq.6. The second term is the posterior of  $\mathbf{r}_t$  given input up to the previous step and the higher-level state:

$$p(\mathbf{r}_t \mid \mathbf{I}_{0:t-1}, \mathbf{r}^h) = \int p(\mathbf{r}_t \mid \mathbf{r}_{t-1}, \mathbf{r}^h) p(\mathbf{r}_{t-1} \mid \mathbf{I}_{0:t-1}, \mathbf{r}^h) d\mathbf{r}_{t-1},$$
(11)

where the first term inside the integral is the lower-level transition dynamics defined by Eq 5. Note that the parameterization of the transition distribution is generated by the higher-level latent variable as specified by Eqs 7 to 9.

Putting Eqs 10 and 11 together, we get

$$p(\mathbf{r}_t \mid \mathbf{I}_{0:t}, \mathbf{r}^h) \propto p(\mathbf{I}_t \mid \mathbf{r}_t) \int p(\mathbf{r}_t \mid \mathbf{r}_{t-1}, \mathbf{r}^h) p(\mathbf{r}_{t-1} \mid \mathbf{I}_{0:t-1}, \mathbf{r}^h) d\mathbf{r}_{t-1},$$
(12)

which defines a recursive way to infer the posterior of  $\mathbf{r}_t$  at time t. In this model, we only maintain a single point (MAP) estimate of the posterior at each time step, so we simplify the posterior distribution  $p(\mathbf{r}_{t-1} \mid \mathbf{I}_{0:t-1}, \mathbf{r}^h)$  as a Dirac delta function:

$$p(\mathbf{r}_{t-1} \mid \mathbf{I}_{0:t-1}, \mathbf{r}^h) \approx \delta(\mathbf{r}_{t-1} - \hat{\mathbf{r}}_{t-1}), \tag{13}$$

where  $\hat{\mathbf{r}}_{t-1}$  is the MAP estimate from the previous step. Now we can further simplify Eq 12 as

$$p(\mathbf{r}_t \mid \mathbf{I}_{0:t}, \mathbf{r}^h) \propto p(\mathbf{I}_t \mid \mathbf{r}_t) p(\mathbf{r}_t \mid \hat{\mathbf{r}}_{t-1}, \mathbf{r}^h). \tag{14}$$

This gives the posterior of all the latent variables at time *t* as

$$p(\mathbf{r}_{t}, \mathbf{r}^{h} \mid \mathbf{I}_{0:t}) \propto p(\mathbf{I}_{t} \mid \mathbf{r}_{t})p(\mathbf{r}_{t} \mid \hat{\mathbf{r}}_{t-1}, \mathbf{r}^{h})p(\mathbf{r}^{h} \mid \mathbf{I}_{0:t}). \tag{15}$$

We can find the MAP estimates of the latent variables by minimizing the negative log of Eq 15. Substituting the generative assumptions (Eqs 4 to 6), we get:

$$\bar{\mathbf{r}}_{t} = \mathbf{f}(\mathbf{r}^{h}, \hat{\mathbf{r}}_{t-1}) \tag{16}$$

$$\mathcal{L}_{t} = \frac{1}{2\sigma^{2}} \|\mathbf{I}_{t} - \mathbf{U}\mathbf{r}_{t}\|_{2}^{2} + \frac{1}{2\sigma_{r}^{2}} \|\mathbf{r}_{t} - \bar{\mathbf{r}}_{t}\|_{2}^{2} + \lambda \|\mathbf{r}_{t}\|_{1} + \lambda_{h} \|\mathbf{r}^{h}\|_{2}^{2},$$
(17)

where  $\lambda$  and  $\lambda_h$  are the sparsity penalty for  $\mathbf{r}_t$  and the Gaussian prior penalty for  $\mathbf{r}^h$ , respectively. Note that we approximate  $p(\mathbf{r}^h \mid \mathbf{I}_{0:t})$  with the unconditional prior  $p(\mathbf{r}^h)$  so that at each step the dynamics are estimated using only the local pairwise transition and the prior. Using Eq 17, we compute the MAP estimate of  $\mathbf{r}_t$  at time t as

$$\hat{\mathbf{r}}_t = \arg\min_{\mathbf{r}_t} \mathcal{L}_t. \tag{18}$$

At each time step, we update the current estimate of  $\mathbf{r}^h$  to minimize  $\mathcal{L}_t$  as well:

$$\mathbf{r}^h = \arg\min_{\mathbf{r}^h} \mathcal{L}_t. \tag{19}$$

To begin the recursive estimation (without the temporal prediction from the previous step), we compute the MAP estimate of the first step latent variable  $\mathbf{r}_0$  using the following reduced loss

$$\mathcal{L}_{0} = \frac{1}{2\sigma^{2}} \|\mathbf{I}_{0} - \mathbf{U}\mathbf{r}_{0}\|_{2}^{2} + \lambda \|\mathbf{r}_{0}\|_{1}$$
(20)

$$\hat{\mathbf{r}}_0 = \arg\min_{\mathbf{r}_0} \mathcal{L}_0. \tag{21}$$

The parameters of the model can be optimized by minimizing the same prediction errors summed across time and averaged across different sequences, using the MAP estimates of the latent variables. See Algorithm A in <u>S1 Text</u> for detailed pseudocode describing the inference and learning procedure.

## Data and preprocessing

For the natural video dataset, we extracted 65520 image sequences from a YouTube video (link here) recorded by a person walking on a forest trail (image size:  $16 \times 16$  pixels, sequence length: 10 frames ( $\approx 0.35$  seconds, uniformly sampled in time). The image sequences do not overlap with each other spatially or temporally. Each sequence was spatially and temporally whitened to simulate retinal and LGN processing following the methods in Olshsausen & Field [38] and Dong & Atick [95]. 58,968 sequences were used to train the model and the remaining 6,552 were reserved for testing.

For the Moving MNIST dataset [45], we used 10000 image sequences (image size:  $18 \times 18$  pixels, sequence length: 10 frames), each sequence containing a fixed digit moving in a particular direction. The motion of the digits was restricted to upward, downward, leftward, or rightward directions. When a digit hit the boundary, its motion direction was inverted (leftward to rightward, upward to downward, and vice versa). No whitening procedures were performed on the MNIST sequences. 9,000 sequences were used to train the model and the remaining 1,000 were reserved for testing.

# Reverse correlation for computing space-time receptive fields

The reverse correlation stimuli for deriving the space-time receptive fields (Fig 2) were extracted from the same natural video data but without any spatial and temporal overlapping with the training and test set. We used 50 continuous image sequences with 80,000 steps ( $\approx$  47 minutes, spanned the same time range but no spatial overlaps) and computed the space-time receptive fields (STRFs) as the firing-rate-weighted average of input frame sequences of length seven ( $\approx$  250 ms), across time and sequences.

Formally, let  $\{\mathbf{I}_{0:T-1}^{(j)}\}_{j=1}^{J}$  be the J stimulus sequences of length T and let  $\tau$  be the length of the STRFs (here, J = 50, T = 80, 000,  $\tau = 7$ ). For each neuron i, its space-time receptive field STRF $_i$  has dimensions  $M \times \tau$ , where M is the dimensionality of a single image frame vector (here, M = 100 after vectorizing the  $10 \times 10$  image frame). We compute the STRF of neuron i as follows

$$STRF_{i} = \frac{1}{J(T-\tau)} \sum_{i=1}^{J} \sum_{t=\tau}^{T-1} \hat{r}_{t,i}^{(j)} \mathbf{I}_{t-\tau:t-1}^{(j)},$$
(22)

where  $\hat{r}_{t,i}^{(j)}$  is the predicted firing rate of neuron i at time t in sequence j and  $\mathbf{I}_{t-\tau:t-1}^{(j)}$  is the image sequence from time  $t-\tau$  to t-1 in sequence j. This procedure is analogous to the calculation of the spike-triggered average widely used in neurophysiology [43]; in this case, we computed the average of input sequences weighted by the activity  $\hat{r}$  caused by the sequence.

## Autocorrelation for quantifying timescales

Since our model responds deterministically to the same sequence (MAP estimation), we cannot follow the exact approach of Murray et al. [6] that relies on across-trial variability to the same stimulus. We computed the autocorrelation of single neuron responses to natural videos and Gaussian white noise sequences. We averaged the single-neuron autocorrelations across the lower or higher level, and trials. Formally, let  $r_{t,i,j}$  be the response of neuron i at time t

during trial j. We computed the autocorrelation with lag k as follows:

$$\mu_{i,j} = \frac{1}{T} \sum_{t=0}^{T-1} r_{t,i,j} \tag{23}$$

$$\sigma_{i,j} = \sqrt{\frac{1}{T} \sum_{t=0}^{T-1} (r_{t,i,j} - \mu_{i,j})^2}$$
 (24)

$$\rho_{i,j}(k) = \frac{\sum_{t=0}^{T-k-1} (r_{t,i,j} - \mu_{i,j}) (r_{t+k,i,j} - \mu_{i,j})}{(T-k)\sigma_{i,j}\sigma_{i,j}} \quad \forall i = 1...N.$$
 (25)

To compute the autocorrelation for an entire population at lag k, we took the average of the autocorrelations across all N neurons and J trials:

$$\rho(k) = \frac{1}{NJ} \sum_{i=1}^{N} \sum_{j=1}^{J} \rho_{ij}(k).$$
 (26)

For the results shown in Fig 3b and 3c, we choose J=500, T=50 and computed the autocorrelation with lag k=0. . . 9 for the lower- and higher-level neurons. White noise pixels were i.i.d. samples from  $\mathcal{N}(0,0.0075)$ . We also computed the autocorrelation for both levels with natural videos selected from the same stimuli used for the reverse correlation analysis (note though with natural videos the stationary mean and variance assumption is less valid). To quantify the timescale of the autocorrelation decay, we fitted an exponential decay function  $\rho_{\rm fit}(k)=a$  exp  $(-k/\tau)+b$  to the autocorrelation data on each level (through Scipy optimize.curve\_fit function), where a, b, and  $\tau$  are fitted parameters of the function and  $\tau$  represents the response timescale following the definition by Murray et al. [6].

## Flash-lag and postdiction simulation

We extracted five-step sequences that have a consistent leftward or rightward motion from the Moving MNIST test set sequences (134 sequences in total, see Fig 5e for an example). To simulate the test conditions used by Eagleman & Sejnowski [26], we either used the first three steps of the sequences to infer a motion (dynamics) estimate  $\hat{\mathbf{r}}^h$  (conditions with initial trajectories), or initialized  $\hat{\mathbf{r}}^h$  to zero vectors (conditions without initial trajectories). Depending on the test condition, the moving object stimulus at t = 3 could move following the original trajectory ("Continuous"), remain at the same location ("Stopped"), move in the reversed trajectory ("Reversed"), or disappear shown an empty frame ("Terminated"), shown in Fig 4a. The stimuli for simulating flashes correspond to the "no initial trajectory, terminate case".

The model's percept of either the moving object or the flashed object at t = 3 was computed as the top-down spatiotemporal prediction of the image after correcting the prediction error at t = 3:

$$\bar{\mathbf{w}} = \mathcal{H}_{\theta}(\hat{\mathbf{r}}^h) \tag{27}$$

$$\bar{\mathbf{V}} = \sum_{k=1}^{K} \bar{w}_k \mathbf{V}_k \tag{28}$$

$$\bar{\mathbf{I}}_{3} = \mathbf{U}(\text{ReLU}(\bar{\mathbf{V}}\hat{\mathbf{r}}_{2})). \tag{29}$$

Here,  $\mathcal{H}_{\theta}$  and  $\mathbf{V}_1, \ldots, \mathbf{V}_K$  are the parameters defined in Eqs 7 and 8,  $\hat{\mathbf{r}}^h$  is the optimal higher-level estimate at t = 3 (Eq 19), and  $\hat{\mathbf{r}}_2$  is the optimal lower-level estimate at t = 2. We computed the location of the percept as the center of mass of the percept image  $\bar{\mathbf{I}}_3$ . The displacement in percept between the moving object and the flashed object was calculated as

$$Displacement = \begin{cases} CoM - x \Big( \bar{\mathbf{I}}_3^{moving} \Big) - CoM - x \Big( \bar{\mathbf{I}}_3^{flash} \Big) & \text{if the moving object has rightward motion} \\ CoM - x \Big( \bar{\mathbf{I}}_3^{flash} \Big) - CoM - x \Big( \bar{\mathbf{I}}_3^{moving} \Big) & \text{if the moving object has leftward motion} \end{cases}$$

where CoM-x(I) returns the horizontal location (x dimension) of the center of mass of I. Therefore, a positive displacement is along the original trajectory of the moving object, while a negative displacement is along the reversed trajectory.

To compute the plots shown in Fig 4g and 4h, we used the "with initial trajectory, reversal" test condition. The displacement was computed as in Eq 30 but  $\bar{\mathbf{I}}_3^{\text{flash}}$  was replaced by the input image  $\mathbf{I}_2$  at t=2, at every value of  $\mathbf{r}^h$  through the error correction process at t=3 (Eqs 27 to 30, see Fig 4b Perceived versus Input).

## Sequence learning and recall simulation

To simulate the sequence learning experiment of Xu et al. [1], we used a five-step sequence extracted from a Moving MNIST test set sequence (Fig 5e). We augmented the hierarchical generative model of DPC with an associative memory layer  $\mathbf{m}$ , which implements predictive coding of the joint higher-level state  $\mathbf{r}^h$  and the lower-level state  $\mathbf{r}_0$  through synapses  $\mathbf{G}$  [10, 58] (see "Summary of the memory model" in S1 Text for model details):

$$(\mathbf{r}_0, \mathbf{r}^h) \mid \mathbf{m} \sim \mathcal{N}(\mathbf{Gm}, \sigma_m^2 I).$$
 (31)

The memory layer was trained separately (the DPC network weights were fixed during conditioning and recall) by minimizing the prediction error:

$$\mathcal{L}_{\text{memory}}(\mathbf{G}, \mathbf{m}) = \|\mathbf{s} - \mathbf{G}\mathbf{m}\|_{2}^{2} + \lambda_{m} \|\mathbf{m}\|_{2}^{2}, \tag{32}$$

where  $\mathbf{s} = [\hat{\mathbf{r}}_0; \hat{\mathbf{r}}^h]$  is the concatenated lower- and higher-level state estimates from DPC and  $\lambda_m$  is the regularization penalty on  $\mathbf{m}$ .

During memory encoding,  $\hat{\mathbf{r}}_0$  and  $\hat{\mathbf{r}}^h$  were estimated by the two-level DPC network from the sequence shown in Fig 5e. Then the memory vector  $\mathbf{m}$  was estimated by gradient descent on Eq 32, yielding the optimal estimate  $\hat{\mathbf{m}}$ :

$$\hat{\boldsymbol{m}} = arg \min_{\boldsymbol{m}} \mathcal{L}_{memory}(\boldsymbol{G}, \boldsymbol{m}). \tag{33}$$

The remaining prediction error drives rapid synaptic plasticity in G through gradient descent on the same equation (Fig 5c):

$$\mathbf{G}' \leftarrow \mathbf{G} - \eta_G \frac{\partial \mathcal{L}_{\text{memory}}(\mathbf{G}, \hat{\mathbf{m}})}{\partial \mathbf{G}},$$
 (34)

where  $\eta_G$  is the learning rate for **G**. During conditioning, we updated the synaptic weights five times using Eq 34. During recall, given the cue for the beginning of the sequence, the full memory vector **m** is retrieved by minimizing the prediction error with respect to the initial lower-level state  $\mathbf{r}_0$  portion of the stored memory only [58] (Fig 5d):

$$\tilde{\mathbf{m}} = \arg\min_{\mathbf{m}} \|\tilde{\mathbf{s}} - \mathbf{b} \odot (\mathbf{G}\mathbf{m})\|_{2}^{2} + \lambda_{m} \|\mathbf{m}\|_{2}^{2}, \tag{35}$$

where  $\tilde{\mathbf{s}} = [\hat{\mathbf{r}}_0; \mathbf{0}^{N_h}]$  denotes the partial input (visual cue representing the first element of the sequence),  $\odot$  denotes element-wise multiplication, and  $\mathbf{b} \in \{0, 1\}^{N+N_h}$  is a binary mask:

$$b_i = \begin{cases} 1 & \text{if } i \le N \\ 0 & \text{otherwise} \end{cases}$$
 (36)

The rest of the sequence was retrieved by retrieving the stored higher-level state  $\bar{\mathbf{r}}^h$  (the dynamics of the sequence) as the top-down prediction through  $\tilde{\mathbf{m}}$  (Fig 5d). Once the dynamics were retrieved, we tested sequential recall in the network by predicting an entire five-step sequence using the lower-level vector  $\hat{\mathbf{r}}_0$  from the visual cue and the retrieved dynamics vector  $\bar{\mathbf{r}}^h$  from the memory (Eqs 7 to 9). We tested recall in the network using three different visual cues: the starting frame (t = 0), the middle frame (t = 0), or the end frame (t = 0) (see Fig 5e and 5f). The cross correlation plot was computed following the same procedure as the one described in Xu et al. [1] (Fig 5g and 5h).

#### Three-level DPC model

To make the level notation clear, we denote the first and second-level states  $\mathbf{r}$  and  $\mathbf{r}^h$  from the two-level model as  $\mathbf{r}^{(1)}$  and  $\mathbf{r}^{(2)}$  in the three-level model, and denote the highest (third-level) state as  $\mathbf{r}^{(3)}$ . We use superscripts to denote level and subscripts to denote time, unless noted otherwise. The trainable parameters for the three-level model include those for the two-level model as well as two second-level transition matrices  $\{\mathbf{V}_1^{(2)},\mathbf{V}_2^{(2)}\}$  and the third-level top-down network  $\mathcal{H}_{\theta}^{(2)}$  that generates the second-level modulation weights.

**Pretraining the second-level transition matrices.** Before training the three-level network, we first pretrained two two-level DPC networks, each with a single transition matrix  $\mathbf{V}^{(2)}$  on Moving MNIST sequences with either straight bouncing type or clockwise bouncing type. We performed inference on the second-level states with the following loss function:

$$\mathcal{L}_{t} = \frac{1}{2\sigma^{2}} \|\mathbf{I}_{t} - \mathbf{U}\mathbf{r}_{t}^{(1)}\|_{2}^{2} + \frac{1}{2\sigma_{r}^{2}} \|\mathbf{r}_{t}^{(1)} - \bar{\mathbf{r}}_{t}^{(1)}\|_{2}^{2} + b_{t} \left(\frac{1}{2\sigma_{r^{(2)}}^{2}} \|\mathbf{r}_{t}^{(2)} - \bar{\mathbf{r}}_{t}^{(2)}\|_{2}^{2}\right) + \lambda \|\mathbf{r}_{t}\|_{1}, \quad (37)$$

where  $\sigma^2_{\mathbf{r}^{(2)}}$  is the variance of second-level prediction errors,  $b_t \in \{0, 1\}$  is a binary mask that equals 1 if the first-level prediction error is larger than a threshold estimated from the training set and 0 otherwise (Fig C in S1 Text), and  $\bar{\mathbf{r}}_t^{(2)} = \mathbf{V}^{(2)}\mathbf{r}_{t-1}^{(2)}$  is the predicted second-level state. We learned  $\mathbf{V}^{(2)}$  by gradient descent on the same loss as in Eq 37, summed across time and averaged across sequences, using the MAP estimates of the latent variables.

**Three-level model training.** We inferred the second- and third-level states using the same loss function (Eq 37), with the predicted second-level state  $\bar{\mathbf{r}}_t^{(2)}$  defined as

$$\mathbf{w}^{(2)} = \mathcal{H}_{\theta}^{(2)}(\mathbf{r}^{(3)}) \tag{38}$$

$$\mathbf{V}^{(2)} = \sum_{k=1}^{K^{(2)}} w_k^{(2)} \mathbf{V}_k^{(2)} \tag{39}$$

$$\bar{\mathbf{r}}_{t}^{(2)} = \mathbf{V}^{(2)} \mathbf{r}_{t-1}^{(2)},$$
 (40)

where  $K^{(2)} = 2$  is the number of second-level transition matrices. Comparing this definition with Eqs  $\sqrt{2}$  to  $\sqrt{2}$ , it is easy to see that the third-level top-down prediction is recursively defined in the same way as the top-down prediction in the two-level model. After obtaining the MAP

estimates of the second- and third-level states, we learned the third-level top-down network  $\mathcal{H}_{\theta}^{(2)}$  by gradient descent on the same loss. See Algorithm B in <u>S1 Text</u> for detailed pseudocode describing the inference and learning procedure for the three-level model.

# **Supporting information**

S1 Text. Supporting information. Fig A. Improvement on test set loss saturates as the **number of transition matrices increases.** (a) Test set loss as training proceeded. Shaded area denotes ±1 standard deviation computed over eight runs with random initialization for each K. K = 1 shows the performance of the single-layer model. (b) Best test loss as K increases. Error bars denote ±1 standard deviation. Fig B. Cue-triggered recall is cue-specific. Four examples of cue-specific sequence recall by the associative memory model after training on different sequences, when given the first frame as the cue. In each quadrant: top: the original image sequence; bottom: cue-triggered recall of the stored sequence. Fig C. Prediction error threshold robustly finds changes of dynamics. (a) The distribution of first-level prediction errors in the two-level DPC model on the Moving MNIST training set. The red dashed line denotes the threshold  $\rho = 0.73$ , where the cumulative density reaches 0.75. (b) Examples of input sequences in the test set. The red arrows mark time steps when the first-level prediction errors exceeded  $\rho$ , corresponding to changes in input dynamics. **Table A. DPC generative** model parameters and values. Table B. Optimizers and learning rates used for inference and learning in the DPC experiments. Here  $\Delta$  denotes the difference in  $\mathbf{r}_t$  or  $\mathbf{r}^h$  from before and after the current iteration of gradient descent. Table C. Memory model parameters and values. Table D. Optimizers and learning rates used for inference and learning in the **memory model**. Here  $\Delta$  denotes the difference in **m** from before and after the current iteration of gradient descent. Table E. Additional parameters and values for the three-level DPC model. Table F. Additional optimizers and learning rates used for inference and learning in the three-level DPC experiments. Algorithm A. Inference & learning process. Algorithm B. Inference & learning process for the three-level DPC model. (PDF)

# **Acknowledgments**

The authors would like to thank Ares Fisher, Dimitrios Gklezakos, Prashant Rangarajan and Vishwas Sathish for discussions related to hypernetworks and predictive coding. LPJ thanks Daogao Liu for inspiring discussions on the modeling aspects of the paper.

#### **Author Contributions**

**Conceptualization:** Linxing Preston Jiang, Rajesh P. N. Rao.

**Data curation:** Linxing Preston Jiang. **Formal analysis:** Linxing Preston Jiang. **Funding acquisition:** Rajesh P. N. Rao.

**Investigation:** Linxing Preston Jiang, Rajesh P. N. Rao. **Methodology:** Linxing Preston Jiang, Rajesh P. N. Rao.

Project administration: Linxing Preston Jiang, Rajesh P. N. Rao.

Resources: Linxing Preston Jiang, Rajesh P. N. Rao.

Software: Linxing Preston Jiang.

**Supervision:** Rajesh P. N. Rao.

Validation: Linxing Preston Jiang.

Visualization: Linxing Preston Jiang.

Writing - original draft: Linxing Preston Jiang, Rajesh P. N. Rao.

Writing – review & editing: Linxing Preston Jiang, Rajesh P. N. Rao.

## References

- Xu S, Jiang W, Poo Mm, Dan Y. Activity recall in a visual cortical ensemble. Nature Neuroscience. 2012; 15(3):449–455. https://doi.org/10.1038/nn.3036 PMID: 22267160
- Keller G, Bonhoeffer T, Hübener M. Sensorimotor Mismatch Signals in Primary Visual Cortex of the Behaving Mouse. Neuron. 2012; 74(5):809–815. <a href="https://doi.org/10.1016/j.neuron.2012.03.040">https://doi.org/10.1016/j.neuron.2012.03.040</a> PMID: 22681686
- Gavornik JP, Bear MF. Learned spatiotemporal sequence recognition and prediction in primary visual cortex. Nature Neuroscience. 2014; 17(5):732–737. https://doi.org/10.1038/nn.3683 PMID: 24657967
- Fiser A, Mahringer D, Oyibo HK, Petersen AV, Leinweber M, Keller GB. Experience-dependent spatial expectations in mouse visual cortex. Nature Neuroscience. 2016; 19(12):1658–1664. https://doi.org/10. 1038/nn.4385 PMID: 27618309
- Schneider DM, Sundararajan J, Mooney R. A cortical filter that learns to suppress the acoustic consequences of movement. Nature. 2018; 561(7723):391–395. <a href="https://doi.org/10.1038/s41586-018-0520-5">https://doi.org/10.1038/s41586-018-0520-5</a> PMID: 30209396
- Murray JD, Bernacchia A, Freedman DJ, Romo R, Wallis JD, Cai X, et al. A hierarchy of intrinsic timescales across primate cortex. Nature Neuroscience. 2014; 17(12):1661–1663. https://doi.org/10.1038/ nn.3862 PMID: 25383900
- Runyan CA, Piasini E, Panzeri S, Harvey CD. Distinct timescales of population coding across cortex. Nature. 2017; 548(7665):92–96. https://doi.org/10.1038/nature23020 PMID: 28723889
- Siegle JH, Jia X, Durand S, Gale S, Bennett C, Graddis N, et al. Survey of spiking in the mouse visual system reveals functional hierarchy. Nature. 2021; 592(7852):86–92. https://doi.org/10.1038/s41586-020-03171-x PMID: 33473216
- Brunec IK, Momennejad I. Predictive Representations in Hippocampal and Prefrontal Hierarchies. Journal of Neuroscience. 2022; 42(2):299–312. https://doi.org/10.1523/JNEUROSCI.1327-21.2021 PMID: 34799416
- Rao RPN, Ballard DH. Predictive coding in the visual cortex: a functional interpretation of some extraclassical receptive-field effects. Nature Neuroscience. 1999; 2(1):79–87. https://doi.org/10.1038/4580 PMID: 10195184
- Friston K. The free-energy principle: a unified brain theory? Nature Reviews Neuroscience. 2010; 11 (2):127–138. https://doi.org/10.1038/nrn2787 PMID: 20068583
- Huang Y, Rao RPN. Predictive coding. WIREs Cognitive Science. 2011; 2(5):580–593. https://doi.org/ 10.1002/wcs.142 PMID: 26302308
- Keller GB, Mrsic-Flogel TD. Predictive Processing: A Canonical Cortical Computation. Neuron. 2018; 100(2):424–435. https://doi.org/10.1016/j.neuron.2018.10.003 PMID: 30359606
- Jiang LP, Rao, Rajesh P N. Predictive Coding Theories of Cortical Function. Oxford Research Encyclopedia of Neuroscience. 2022;.
- **15.** Ha D, Dai AM, Le QV. HyperNetworks. In: 5th International Conference on Learning Representations (ICLR 2017); 2017.
- Ferguson KA, Cardin JA. Mechanisms underlying gain modulation in the cortex. Nature Reviews Neuroscience. 2020; 21(2):80–92. https://doi.org/10.1038/s41583-019-0253-y PMID: 31911627
- Shine JM, Müller EJ, Munn B, Cabral J, Moran RJ, Breakspear M. Computational models link cellular mechanisms of neuromodulation to large-scale neural dynamics. Nature Neuroscience. 2021; 24 (6):765–776. https://doi.org/10.1038/s41593-021-00824-6 PMID: 33958801
- van Hateren JH, Ruderman DL. Independent component analysis of natural image sequences yields spatio-temporal filters similar to simple cells in primary visual cortex. Proceedings of the Royal Society of London Series B: Biological Sciences. 1998; 265:2315–2320. https://doi.org/10.1098/rspb.1998. 0577 PMID: 9881476

- Kayser C, Einhäuser W, Dümmer O, König P, Körding K. Extracting Slow Subspaces from Natural Videos Leads to Complex Cells. In: International Conference on Artificial Neural Networks. Lecture Notes in Computer Science; 2001. p. 1075–1080.
- Olshausen BA. Sparse coding of time-varying natural images. Journal of Vision. 2002; 2(7):130–130. https://doi.org/10.1167/2.7.130
- Wiskott L, Sejnowski TJ. Slow Feature Analysis: Unsupervised Learning of Invariances. Neural Computation. 2002; 14(4):715–770. https://doi.org/10.1162/089976602317318938 PMID: 11936959
- 22. Berkes P, Wiskott L. Slow feature analysis yields a rich repertoire of complex cell properties. Journal of Vision. 2005; 5(6):9. https://doi.org/10.1167/5.6.9 PMID: 16097870
- Singer Y, Teramoto Y, Willmore BD, Schnupp JW, King AJ, Harper NS. Sensory cortex is optimized for prediction of future input. eLife. 2018; 7:e31557. https://doi.org/10.7554/eLife.31557 PMID: 29911971
- 24. Singer Y, Willmore BDB, King AJ, Harper NS. Hierarchical temporal prediction captures motion processing from retina to higher visual cortex; 2019. Available from: http://biorxiv.org/lookup/doi/10.1101/575464.
- DeAngelis GC, Ohzawa I, Freeman RD. Receptive-field dynamics in the central visual pathways. Trends in Neurosciences. 1995; 18(10):451–458. https://doi.org/10.1016/0166-2236(95)94496-R PMID: 8545912
- Eagleman DM, Sejnowski TJ. Motion Integration and Postdiction in Visual Awareness. Science. 2000; 287(5460):2036–2038. https://doi.org/10.1126/science.287.5460.2036 PMID: 10720334
- Hogendoorn H, Carlson TA, Verstraten FAJ. Interpolation and extrapolation on the path of apparent motion. Vision Research. 2008; 48(7):872–881. <a href="https://doi.org/10.1016/j.visres.2007.12.019">https://doi.org/10.1016/j.visres.2007.12.019</a> PMID: 18279906
- Shimojo S. Postdiction: its implications on visual awareness, hindsight, and sense of agency. Frontiers in Psychology. 2014; 5. https://doi.org/10.3389/fpsyg.2014.00196 PMID: 24744739
- Hogendoorn H. Perception in real-time: predicting the present, reconstructing the past. Trends in Cognitive Sciences. 2022; 26(2):128–141. https://doi.org/10.1016/j.tics.2021.11.003 PMID: 34973925
- Nijhawan R. Motion extrapolation in catching. Nature. 1994; 370(6487):256–257. https://doi.org/10. 1038/370256b0 PMID: 8035873
- Nijhawan R. Visual prediction: Psychophysics and neurophysiology of compensation for time delays. Behavioral and Brain Sciences. 2008; 31(2):179–198. <a href="https://doi.org/10.1017/S0140525X08003804">https://doi.org/10.1017/S0140525X08003804</a> PMID: 18479557
- **32.** Ekman M, Kok P, de Lange FP. Time-compressed preplay of anticipated events in human primary visual cortex. Nature Communications. 2017; 8(1):15276. https://doi.org/10.1038/ncomms15276 PMID: 28534870
- Bang JW, Sasaki Y, Watanabe T, Rahnev D. Feature-Specific Awake Reactivation in Human V1 after Visual Training. Journal of Neuroscience. 2018; 38(45):9648–9657. <a href="https://doi.org/10.1523/JNEUROSCI.0884-18.2018">https://doi.org/10.1523/JNEUROSCI.0884-18.2018</a> PMID: 30242054
- Lu J, Luo L, Wang Q, Fang F, Chen N. Cue-triggered activity replay in human early visual cortex. Science China Life Sciences. 2021; 64(1):144–151. <a href="https://doi.org/10.1007/s11427-020-1726-5">https://doi.org/10.1007/s11427-020-1726-5</a> PMID: 32557289
- Eagleman SL, Dragoi V. Image sequence reactivation in awake V4 networks. Proceedings of the National Academy of Sciences. 2012; 109(47):19450–19455. https://doi.org/10.1073/pnas. 1212059109 PMID: 23129638
- **36.** Jiang LP, Rao RPN. Dynamic Predictive Coding Explains Both Prediction and Postdiction in Visual Motion Perception. Proceedings of the Annual Meeting of the Cognitive Science Society. 2023;45(45).
- George D, Hawkins J. Towards a Mathematical Theory of Cortical Micro-circuits. PLOS Computational Biology. 2009; 5(10):e1000532. https://doi.org/10.1371/journal.pcbi.1000532 PMID: 19816557
- Olshausen BA, Field DJ. Emergence of simple-cell receptive field properties by learning a sparse code for natural images. Nature. 1996; 381(6583):607–609. <a href="https://doi.org/10.1038/381607a0">https://doi.org/10.1038/381607a0</a> PMID: 8637596
- Larkum ME, Senn W, Lüscher HR. Top-down Dendritic Input Increases the Gain of Layer 5 Pyramidal Neurons. Cerebral Cortex. 2004; 14(10):1059–1070. <a href="https://doi.org/10.1093/cercor/bhh065">https://doi.org/10.1093/cercor/bhh065</a> PMID: 15115747
- 40. Rao RPN. An optimal estimation approach to visual perception and learning. Vision Research. 1999; 39 (11):1963–1989. https://doi.org/10.1016/S0042-6989(98)00279-X PMID: 10343783
- Dong DW, Atick JJ. Statistics of natural time-varying images. Network: Computation in Neural Systems. 1995; 6(3):345–358. https://doi.org/10.1088/0954-898X\_6\_3\_003
- Hubel DH, Wiesel TN. Receptive fields of single neurones in the cat's striate cortex. The Journal of Physiology. 1959; 148(3):574–591. https://doi.org/10.1113/jphysiol.1959.sp006308 PMID: 14403679

- Ringach D, Shapley R. Reverse correlation in neurophysiology. Cognitive Science. 2004; 28(2):147– 166. https://doi.org/10.1207/s15516709cog2802\_2
- Talebi V, Baker CL. Natural versus Synthetic Stimuli for Estimating Receptive Field Models: A Comparison of Predictive Robustness. Journal of Neuroscience. 2012; 32(5):1560–1576. <a href="https://doi.org/10.1523/JNEUROSCI.4661-12.2012">https://doi.org/10.1523/JNEUROSCI.4661-12.2012</a> PMID: 22302799
- 45. Srivastava N, Mansimov E, Salakhudinov R. Unsupervised Learning of Video Representations using LSTMs. In: Proceedings of the 32nd International Conference on Machine Learning; 2015. p. 843–852. Available from: https://proceedings.mlr.press/v37/srivastava15.html.
- 46. Murphy KP. Machine learning: a probabilistic perspective. Cambridge, Massachusetts: MIT Press; 2013. Available from: https://www.amazon.com/Machine-Learning-Probabilistic-Perspective-Computation/dp/0262018020/ref=sr\_1\_2?ie=UTF8&qid=1336857747&sr=8-2.
- 47. Jordan R, Keller GB. Opposing Influence of Top-down and Bottom-up Input on Excitatory Layer 2/3 Neurons in Mouse Primary Visual Cortex. Neuron. 2020; 108(6):1194–1206.e5. <a href="https://doi.org/10.1016/j.neuron.2020.09.024">https://doi.org/10.1016/j.neuron.2020.09.024</a> PMID: 33091338
- **48.** Ekman M, Gennari G, Lange FPd. Probabilistic forward replay of anticipated stimulus sequences in human primary visual cortex and hippocampus; 2022. Available from: <a href="https://www.biorxiv.org/content/10.1101/2022.01.26.477907v1">https://www.biorxiv.org/content/10.1101/2022.01.26.477907v1</a>.
- **49.** Rao RPN. Correlates of Attention in a Model of Dynamic Visual Recognition. In: Advances in Neural Information Processing Systems; 1998. Available from: <a href="http://papers.nips.cc/paper/1416-correlates-of-attention-in-a-model-of-dynamic-visual-recognition.pdf">http://papers.nips.cc/paper/1416-correlates-of-attention-in-a-model-of-dynamic-visual-recognition.pdf</a>.
- Manns JR, Eichenbaum H. Evolution of declarative memory. Hippocampus. 2006; 16(9):795–808. https://doi.org/10.1002/hipo.20205 PMID: 16881079
- Burgess N, Maguire EA, O'Keefe J. The Human Hippocampus and Spatial and Episodic Memory. Neuron. 2002; 35(4):625–641. https://doi.org/10.1016/S0896-6273(02)00830-9 PMID: 12194864
- Tulving E. Episodic memory: From mind to brain. Annual Review of Psychology. 2002; 53(1):1–25. https://doi.org/10.1146/annurev.psych.53.100901.135114 PMID: 11752477
- Sugar J, Moser MB. Episodic memory: Neuronal codes for what, where, and when. Hippocampus. 2019; 29(12):1190–1205. https://doi.org/10.1002/hipo.23132 PMID: 31334573
- 54. Gelbard-Sagiv H, Mukamel R, Harel M, Malach R, Fried I. Internally Generated Reactivation of Single Neurons in Human Hippocampus During Free Recall. Science. 2008; 322(5898):96–101. <a href="https://doi.org/10.1126/science.1164685">https://doi.org/10.1126/science.1164685</a> PMID: 18772395
- 55. Bosch SE, Jehee JFM, Fernández G, Doeller CF. Reinstatement of Associative Memories in Early Visual Cortex Is Signaled by the Hippocampus. Journal of Neuroscience. 2014; 34(22):7493–7500. https://doi.org/10.1523/JNEUROSCI.0805-14.2014 PMID: 24872554
- Hindy NC, Ng FY, Turk-Browne NB. Linking pattern completion in the hippocampus to predictive coding in visual cortex. Nature Neuroscience. 2016; 19(5):665–667. https://doi.org/10.1038/nn.4284 PMID: 27065363
- Barron HC, Auksztulewicz R, Friston K. Prediction and memory: A predictive coding account. Progress in Neurobiology. 2020; 192:101821. <a href="https://doi.org/10.1016/j.pneurobio.2020.101821">https://doi.org/10.1016/j.pneurobio.2020.101821</a> PMID: 32446883
- Salvatori T, Song Y, Hong Y, Sha L, Frieder S, Xu Z, et al. Associative Memories via Predictive Coding. In: Advances in Neural Information Processing Systems. vol. 34; 2021. p. 3874–3886.
- 59. Fine S, Singer Y, Tishby N. The Hierarchical Hidden Markov Model: Analysis and Applications. Machine Learning. 1998; 32(1):41–62. https://doi.org/10.1023/A:1007469218079
- Hogendoorn H. Motion Extrapolation in Visual Processing: Lessons from 25 Years of Flash-Lag Debate. Journal of Neuroscience. 2020; 40(30):5698–5705. https://doi.org/10.1523/JNEUROSCI.0275-20.2020 PMID: 32699152
- Mishkin M, Ungerleider LG, Macko KA. Object vision and spatial vision: two cortical pathways. Trends in Neurosciences. 1983; 6:414–417. https://doi.org/10.1016/0166-2236(83)90190-X
- Mante V, Sussillo D, Shenoy KV, Newsome WT. Context-dependent computation by recurrent dynamics in prefrontal cortex. Nature. 2013; 503(7474):78–84. <a href="https://doi.org/10.1038/nature12742">https://doi.org/10.1038/nature12742</a> PMID: 24201281
- Stroud JP, Porter MA, Hennequin G, Vogels TP. Motor primitives in space and time via targeted gain modulation in cortical networks. Nature Neuroscience. 2018; 21(12):1774–1783. <a href="https://doi.org/10.1038/s41593-018-0276-0">https://doi.org/10.1038/s41593-018-0276-0</a> PMID: 30482949
- 64. Masse NY, Grant GD, Freedman DJ. Alleviating catastrophic forgetting using context-dependent gating and synaptic stabilization. Proceedings of the National Academy of Sciences. 2018; 115(44):E10467– E10475. https://doi.org/10.1073/pnas.1803839115 PMID: 30315147

- 65. Shai AS, Anastassiou CA, Larkum ME, Koch C. Physiology of Layer 5 Pyramidal Neurons in Mouse Primary Visual Cortex: Coincidence Detection through Bursting. PLOS Computational Biology. 2015; 11 (3):e1004090. https://doi.org/10.1371/journal.pcbi.1004090 PMID: 25768881
- 66. Takahashi N, Oertner TG, Hegemann P, Larkum ME. Active cortical dendrites modulate perception. Science. 2016; 354(6319):1587–1590. https://doi.org/10.1126/science.aah6066 PMID: 28008068
- Galanti T, Wolf L. On the Modularity of Hypernetworks. In: Advances in Neural Information Processing Systems. vol. 33. Curran Associates, Inc.; 2020. p. 10409–10419.
- Srinivasan MV, Laughlin SB, Dubs A. Predictive coding: a fresh view of inhibition in the retina. Proceedings of the Royal Society of London Series B Biological Sciences. 1982; 216(1205):427–459. PMID: 6129637
- Luczak A, McNaughton BL, Kubo Y. Neurons learn by predicting future activity. Nature Machine Intelligence. 2022; 4(1):62–72. https://doi.org/10.1038/s42256-021-00430-y PMID: 35814496
- LeCun Y, Bengio Y, Hinton G. Deep learning. Nature. 2015; 521(7553):436–444. https://doi.org/10. 1038/nature14539 PMID: 26017442
- 71. Lotter W, Kreiman G, Cox DD. Deep predictive coding networks for video prediction and unsupervised learning. In: International Conference on Learning Representations; 2017. Available from: <a href="https://openreview.net/forum?id=B1ewdt9xe">https://openreview.net/forum?id=B1ewdt9xe</a>.
- 72. Lotter W, Kreiman G, Cox D. A neural network trained for prediction mimics diverse features of biological neurons and perception. Nature Machine Intelligence. 2020; 2(4):210–219. https://doi.org/10.1038/s42256-020-0170-9 PMID: 34291193
- 73. Piasini E, Soltuzu L, Muratore P, Caramellino R, Vinken K, Op de Beeck H, et al. Temporal stability of stimulus representation increases along rodent visual cortical hierarchies. Nature Communications. 2021; 12(1):4448. https://doi.org/10.1038/s41467-021-24456-3 PMID: 34290247
- Henin S, Turk-Browne NB, Friedman D, Liu A, Dugan P, Flinker A, et al. Learning hierarchical sequence representations across human cortex and hippocampus. Science Advances. 2021; 7. <a href="https://doi.org/10.1126/sciadv.abc4530">https://doi.org/10.1126/sciadv.abc4530</a> PMID: 33608265
- 75. Chaudhuri R, Knoblauch K, Gariel MA, Kennedy H, Wang XJ. A Large-Scale Circuit Mechanism for Hierarchical Dynamical Processing in the Primate Cortex. Neuron. 2015; 88(2):419–431. https://doi. org/10.1016/j.neuron.2015.09.008 PMID: 26439530
- Joglekar MR, Mejias JF, Yang GR, Wang XJ. Inter-areal Balanced Amplification Enhances Signal Propagation in a Large-Scale Circuit Model of the Primate Cortex. Neuron. 2018; 98(1):222–234.e8. https://doi.org/10.1016/j.neuron.2018.02.031 PMID: 29576389
- van Meegen A, van Albada SJ. Microscopic theory of intrinsic timescales in spiking neural networks. Physical Review Research. 2021; 3(4):043077. https://doi.org/10.1103/PhysRevResearch.3.043077
- Kiebel SJ, Daunizeau J, Friston KJ. A Hierarchy of Time-Scales and the Brain. PLOS Computational Biology. 2008; 4(11):e1000209. https://doi.org/10.1371/journal.pcbi.1000209 PMID: 19008936
- 79. Gao R, van den Brink RL, Pfeffer T, Voytek B. Neuronal timescales are functionally dynamic and shaped by cortical microarchitecture. eLife. 2020; 9:e61277. https://doi.org/10.7554/eLife.61277 PMID: 33226336
- Rao RPN, Eagleman DM, Sejnowski TJ. Optimal Smoothing in Visual Motion Perception. Neural Computation. 2001; 13(6):1243–1253. https://doi.org/10.1162/08997660152002843 PMID: 11387045
- Khoei MA, Masson GS, Perrinet LU. The Flash-Lag Effect as a Motion-Based Predictive Shift. PLOS Computational Biology. 2017; 13(1):e1005068. <a href="https://doi.org/10.1371/journal.pcbi.1005068">https://doi.org/10.1371/journal.pcbi.1005068</a> PMID: 28125585
- **82.** Finnie PSB, Komorowski RW, Bear MF. The spatiotemporal organization of experience dictates hippocampal involvement in primary visual cortical plasticity. Current Biology. 2021; 31(18):3996–4008.e6. https://doi.org/10.1016/j.cub.2021.06.079 PMID: 34314678
- 83. Foster DJ. Replay Comes of Age. Annual Review of Neuroscience. 2017; 40(1):581–602. https://doi.org/10.1146/annurev-neuro-072116-031538 PMID: 28772098
- 84. Friston K. A theory of cortical responses. Philosophical Transactions of the Royal Society B: Biological Sciences. 2005; 360(1456):815–836. https://doi.org/10.1098/rstb.2005.1622 PMID: 15937014
- 85. Linderman S, Johnson M, Miller A, Adams R, Blei D, Paninski L. Bayesian Learning and Inference in Recurrent Switching Linear Dynamical Systems. In: Proceedings of the 20th International Conference on Artificial Intelligence and Statistics. PMLR; 2017. p. 914–922. Available from: https://proceedings. mlr.press/v54/linderman17a.html.
- **86.** Rao RPN, Gklezakos DC, Sathish V. Active Predictive Coding: A Unifying Neural Model for Active Perception, Compositional Learning, and Hierarchical Planning. Neural Computation. 2024; 36(1):1–32. https://doi.org/10.1162/neco\_a\_01627

- Zmarz P, Keller GB. Mismatch Receptive Fields in Mouse Visual Cortex. Neuron. 2016; 92(4):766–772. https://doi.org/10.1016/j.neuron.2016.09.057 PMID: 27974161
- **88.** Attias H. Planning by Probabilistic Inference. In: International Workshop on Artificial Intelligence and Statistics; 2003. p. 9–16. Available from: http://proceedings.mlr.press/r4/attias03a.html.
- **89.** Verma D, Rao RP. Goal-Based Imitation as Probabilistic Inference over Graphical Models. Advances in Neural Information Processing Systems. 2005;18.
- **90.** Verma D, Rao RPN. Planning and Acting in Uncertain Environments using Probabilistic Inference. In: 2006 IEEE/RSJ International Conference on Intelligent Robots and Systems; 2006. p. 2382–2387.
- 91. Botvinick M, Toussaint M. Planning as inference. Trends in Cognitive Sciences. 2012; 16(10):485–488. https://doi.org/10.1016/j.tics.2012.08.006 PMID: 22940577
- Levine S. Reinforcement Learning and Control as Probabilistic Inference: Tutorial and Review. arXiv:180500909 [cs, stat]. 2018;.
- Momennejad I. Learning Structures: Predictive Representations, Replay, and Generalization. Current Opinion in Behavioral Sciences. 2020; 32:155–166. <a href="https://doi.org/10.1016/j.cobeha.2020.02.017">https://doi.org/10.1016/j.cobeha.2020.02.017</a> PMID: 35419465
- Stachenfeld KL, Botvinick MM, Gershman SJ. The hippocampus as a predictive map. Nature Neuroscience. 2017; 20(11):1643–1653. https://doi.org/10.1038/nn.4650 PMID: 28967910
- **95.** Dong DW, Atick JJ. Temporal decorrelation: a theory of lagged and nonlagged responses in the lateral geniculate nucleus. Network: Computation in neural systems. 1995; 6(2):159–178. https://doi.org/10. 1088/0954-898X\_6\_2\_003