

## Reviewed Preprint

v1 • May 18, 2026

Not revised

## ✉ For correspondence:

friedemann.zenke@fmi.ch

\* These authors contributed equally to this work.

## Competing interests: No

competing interests declared

Funding: See [page 28](#)

## Reviewing editor: Tatyana O

Sharpee, Salk Institute for Biological Studies, United States

© 2026, Mohammadi et al. This article is distributed under the terms of the [Creative Commons Attribution License](#), which permits unrestricted use and redistribution provided that the original author and source are credited.

# Understanding neural circuit principles for representation learning through joint-embedding predictive architectures

Ashena Gorgan Mohammadi<sup>1,2,\*</sup>, Manu Srinath Halvagal<sup>1,2,3,\*</sup>, Friedemann Zenke<sup>1,2</sup> ✉

<sup>1</sup>Friedrich Miescher Institute for Biomedical Research, Basel, Switzerland • <sup>2</sup>Faculty of Science, University of Basel, Basel, Switzerland • <sup>3</sup>Department of Physics, Harvard University, Cambridge, United States

## eLife Assessment

This manuscript proposes a **valuable** idea on how cortical networks may learn a helpful representation of sensory stimuli. The model implementing this idea is tested in multiple experimental paradigms. However, the evidence remains **incomplete** as to whether the method supports both invariance and equivariance and whether it can estimate the dynamics of the moving object.

<https://doi.org/10.7554/eLife.110895.1.sa4>

## Abstract

Tracking prey or recognizing a lurking predator is as crucial for survival as anticipating their actions. To guide behavior, the brain must extract information about object identities and their dynamics from entangled sensory inputs. How it accomplishes this feat remains an open question. Classical predictive coding theories propose that this ability arises by comparing predicted sensory signals with actual inputs and reducing the associated prediction errors. While such models capture important aspects of cortical computation, they typically focus on faithfully predicting sensory input and do not explicitly address how abstract, untangled representations of objects and their dynamics emerge solely through experience. Here, we develop a theory of representation learning in neural circuits that shifts the focus from prediction in the input space to prediction in representation space, without relying on external supervision or labeled data. Specifically, we introduce recurrent predictive learning (RPL), a recurrent joint-embedding predictive architecture inspired by self-supervised machine learning, that learns abstract representations of object identity and their dynamics and predicts future object motion from continuous sensory streams. Crucially, the model learns sequence representations that resemble successor-like representations observed in the primary visual cortex of humans. The model also develops abstract sequence representations comparable to those reported in the macaque prefrontal cortex. Finally, we outline how RPL's modular feedforward-recurrent organization could map onto cortical microcircuits. Our work establishes a circuit-centric theory framework that provides new perspectives on how the brain may acquire an internal model of the world through experience.

## Introduction

It takes us a split second to recognize an oncoming car at an intersection, to anticipate the trajectory of a ball thrown to us, or to realize that we are being chased by a dog in the park. Importantly, we have no difficulty assessing these situations and responding to them, despite never having seen this particular car, ball, or dog. To accomplish this feat, our nervous system faces a twofold challenge. On the one hand, it must detect abstract patterns from sensory input and form invariant object representations, allowing us to recognize, for instance, a dog. This task

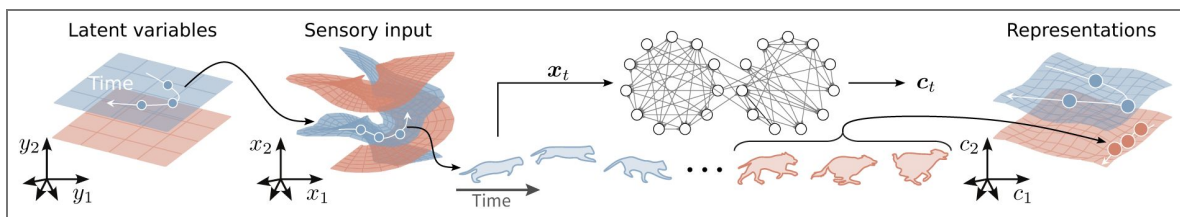
is hindered by the fact that distinct sensory stimuli are often entangled, requiring complex hierarchical networks to untangle them (Fig. 1). On the other hand, our brain must also infer equivariant features of the dog's dynamics, such as its position, velocity, and angle of attack. After all, our actions depend on whether it is moving toward or away from us. While a plethora of experiments have shown that sensory cortices encode this information [1–9], the mechanisms underlying the plasticity that forms these representations remain elusive. Previous modeling work indicated that minimal extensions to local Hebbian learning rules can learn invariant object representations in deep neural networks [10, 11]. However, it remains unclear whether similar mechanisms also yield equivariant representations. Yet, both invariant *and* equivariant representations are necessary to respond appropriately in novel situations.

Another influential idea of how representation learning works in the cortex is predictive processing [12], typified by classical predictive coding models [13]. These propose that the brain acquires an internal model of the world that enables it to anticipate incoming sensory input. In this framework, cortical networks generate predictions of expected sensory stimuli and compare them to the actual input, yielding prediction errors. Minimizing these errors supports perceptual inference of latent variables, representation learning, and the continual refinement of the internal model [13–22].

While a compelling idea, classical predictive coding models have one major drawback: they are generative models that construct predictions in the input space. Consequently, such models require a decoder network that predicts, for example, the pixel-level image of an incoming dog, which is then compared to the actual sensory experience, rather than comparing internal representations directly. This subtle distinction has profound consequences for the neural circuit architecture required, specifically where in the circuit prediction errors are computed, how they interact with representation neurons, and how they guide learning.

A growing body of work has shown that classical predictive coding models trained as fully generative systems often struggle to yield representations that support abstraction, invariance, or object-level discrimination. In particular, models that minimize local prediction errors greedily at each hierarchical stage tend to learn representations that remain closely tied to input statistics, limiting their usefulness for downstream classification or invariant recognition [23, 24]. This limitation is also evident in large-scale models such as PredNet [15]. While PredNet achieves its best performance when prediction errors are optimized only at the input (pixel) level through end-to-end training, thereby allowing deeper representations to be shaped indirectly, its representational and predictive performance degrades when prediction errors are minimized across all hierarchical levels. Together, these findings suggest a tension between generative, reconstruction-driven objectives and the learning of abstract representations, motivating predictive frameworks that operate directly in the representation space rather than the sensory domain.

Finally, classical predictive coding posits that predicted sensory information is represented with reduced neural activity in early sensory cortices due to suppressive top-down feedback. However, a growing body of neurophysiological evidence indicates that top-down input from higher cortical areas can enhance representations of behaviorally relevant or attended stimuli [25]. Consequently, whether and how we can reconcile predictive processing with cortical anatomy, abstract representation learning, and empirical observations of inter-areal interactions remains a central outstanding question in circuit neuroscience [12, 25–30].



**Figure 1. Extracting latent variables from sensory stimuli.**

Animals and humans must infer latent variables  $\mathbf{y}$  such as object identity, orientation, and velocity from ambiguous and noisy sensory stimuli. This task is computationally demanding because sensory stimuli  $\mathbf{x}$  are typically entangled (left). Untangling these input manifolds to approximate the underlying latent variables  $\mathbf{c}_t$  (right) requires multi-layer networks (middle). How our brains achieve this untangling remains a central question in neuroscience.

A natural alternative is to abandon explicit generative modeling and instead compute prediction errors entirely within internal representation space. This strategy has proven effective in recent self-supervised learning frameworks [31]. By eliminating the need for a decoder that reconstructs sensory input, such an approach would avoid the difficulties associated with generative prediction in the input space. However, representation space prediction introduces a new challenge: *representational collapse*, wherein networks converge to trivial solutions, e.g., constant representations, that are perfectly predictable but devoid of useful information. Preventing collapse requires additional constraints, such as variance regularization techniques [11, 32, 33], negative samples [10, 34], or a joint-embedding predictive architecture (JEPA) [35].

JEPAs avoid collapse by computing prediction errors between internal representations using a dedicated predictor network [36–41]. Beyond biological plausibility, this strategy offers another distinct advantage: by dispensing with reconstruction objectives, the network can devote its coding capacity to abstract, behaviorally relevant variables such as object identity and future dynamics, e.g., “the dog behind me is about to pounce,” rather than low-level sensory details that may be incidental or task-irrelevant, such as the particular sheen of the dog’s fur around the ears [42, 43].

Despite these advantages, most JEPAs developed in machine learning are poor models of cortical computation. They rely on architectural elements with limited biological grounding, such as transformers, and training strategies, such as token masking, while focusing mainly on static images, hand-crafted transformations, and invariant representations geared towards object classification. As a result, it remains unclear whether a JEPA broadly consistent with cortical circuit structure and function exists. A minimal candidate must instead rely on recurrent neural dynamics, learn from streaming sensory input without masking, support both invariant and equivariant representations, and reproduce key neurophysiological observations.

Here, we introduce such a model. Our goal is to develop a circuit-level theory of representation learning that explains how abstract representations could emerge in cortical networks and to identify qualitative correspondences and testable predictions, rather than to quantitatively fit neural data. To that end, we propose recurrent predictive learning (RPL). This minimal circuit model learns linearly decodable object representations, including their associated motion variables, from videos of moving objects without requiring any labels. As we will see, it yields abstract representations spanning multiple timescales, consistent with representations in the macaque prefrontal cortex (PFC) [44]. In addition to representation learning, RPL learns a dynamic world model, which represents future stimuli, similar to human V1 [45]. To accomplish this, RPL requires distinct feedforward and recurrent components, each playing a dedicated functional role. In a hierarchical setting, prediction errors are computed in local prediction circuits; this organization allows for an interpretation in terms of necessary neuronal circuit elements. In summary, we introduce a circuit-centric network model and learning strategy that extends the notion of predictive processing to internal representation prediction, providing a principled set of components for learning structured representations in the brain.

## Results

To study the putative mechanisms by which sensory networks extract latent variables, such as invariant object identity, and equivariant latent variables, such as position, velocity, and orientation, we simulated deep neural network models with a specific structure. Broadly speaking, these networks consist of three functionally distinct components (Fig. 2a [↗](#)). First, a feedforward “encoder” network nonlinearly encodes sensory stimuli into the intermediate embedding  $z$  with minimal to no history dependence. Second, a recurrent neural network, the “integrator,” successively integrates these embeddings to form the internal representations  $c$ . Third, the predictor network  $f$  takes the representation  $c_t$  at time  $t$  to predict the embedding  $z_{t+\Delta}$  an instant later. Finally, the resulting prediction error

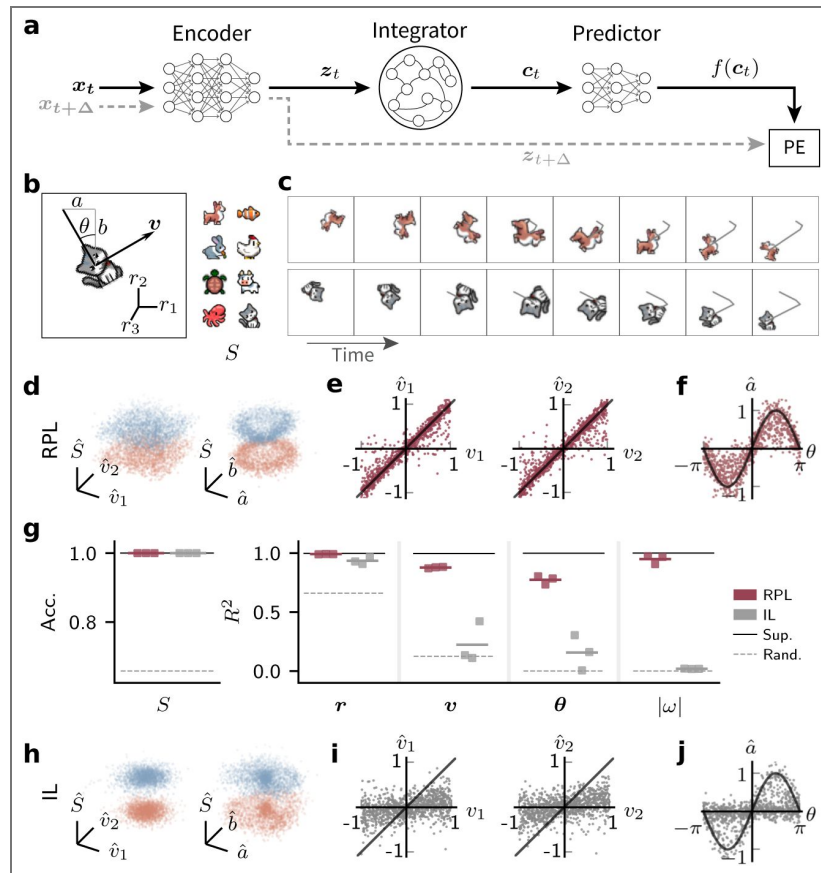
$$e_{t+\Delta} = \underbrace{z_{t+\Delta}}_{\text{embedding}} - \underbrace{f(c_t)}_{\text{prediction}} \quad (1)$$

serves as an instructive learning signal to optimize the model parameters (Methods). Importantly, this unsupervised learning signal is fully specified in the internal neuronal representation space, localized in time, and can be computed without an explicit prediction or reconstruction of  $x$  in the input space. We collectively refer to network models trained by minimizing the above prediction error as RPL in the following.

To evaluate RPL's effectiveness in forming untangled representations of both invariant and dynamic variables from object motion sequences, we constructed a synthetic dataset of moving animal images (Fig. 2b; Methods) with *known* latent variables. To that end, we generated short sequences of one of eight cartoon animal images moving smoothly with a randomly chosen linear and angular velocity within a square arena with elastic collisions at the boundaries (Fig. 2c). To obtain individual video frames, we sampled the trajectory at discrete intervals, placed each frame in front of a gray background, and added Gaussian pixel noise (Supplementary Fig. S1; Methods). Each video frame was characterized by the animal's identity  $S$ , its position  $r$ , and velocity  $v$ , with  $r_{1/2}$  defining the planar position and  $r_3$  defining the scale, orientation  $\theta$ , and the corresponding angular velocity  $\omega$ . Finally, to ensure that these latent variables were indeed entangled in stimulus space, we measured how well a linear decoder could predict them from pixel values alone. To that end, we quantified linear classification accuracy and decoding quality ( $R^2$ ) of dynamic variables. Most variables could not reliably be decoded from the raw input image sequences (Supplementary Fig. S2; Table S1).

To test whether RPL recovers the underlying invariant and dynamic latent variables, we trained our model on approximately 5.7 h of synthetic videos (Methods). We then visualized the learned representations by plotting two example classes (Cat and Dog) onto their best linear projections of the latent linear velocities  $v_{1/2}$ . The resulting point clouds faithfully recovered the underlying uniform velocity distribution of the data (Fig. 2d,e; Supplementary Fig. S3a,b). Similarly, when regressing on  $a = \sin \theta$  and  $b = \cos \theta$ , the resulting linear projections revealed an explicit orientation encoding in the internal representations (Fig. 2d,f). Finally, we quantified classification accuracy and decoding quality of position, velocity, orientation  $\theta = (a, b)$ , and angular speed  $|\omega|$ , and saw high decoding performance for all latent variables (Fig. 2g). The decoding performance of RPL was substantially better than a random untrained baseline model and close to a reference model trained using supervised learning. For comparison to contrastive self-supervised learning (SSL), a powerful machine learning approach that typically violates temporal locality, we also trained and evaluated the same network model using the contrastive predictive coding objective [34]. RPL training yielded better linear decodability of equivariant features (Table S1; Methods). Thus, RPL learns untangled representations of object identity and motion variables that can be read out linearly, a computation that downstream neurons can easily implement.

We were curious whether this dual untangling of object identity and motion variables depends on the particular type of JEPA we used, in which the recurrent pathway predicts the feed-forward signal. To address this question, we kept the network architecture unchanged and used invariance learning (IL) with a symmetric learning objective similar to previous work [11, 32] in which the predictor network predicts its own output instead of the encoder output (Methods). Concretely, this IL objective minimizes the prediction error  $e_{t+\Delta}^{\text{IL}} = f(c_{t+\Delta}) - f(c_t)$ . Networks trained to minimize such IL prediction error formed untangled representations of object identity and position. Still, their representations did not allow reliable decoding of other latent variables such as velocity, orientation, and angular speed (Fig. 2g-j; Supplementary Figs. S3c,d, S4, S5; Supplementary Table S1). Similarly, optimizing the network using an IL objective at the level of the embedding  $z$  or the representation  $c$  (Eqs. (4), (5)), resulted in worse decoding performance than the RPL objective (Supplementary Fig. S2). Thus RPL's asymmetric recurrent prediction circuit is essential to learn representations of static and dynamic latent variables and specifically for representing direction of motion or sequence order.



**Figure 2. Learning latent variables from motion stimuli without reconstruction and labels.**

(a) Schematic of the network model. RPL consists of different network components. First, the “encoder” encodes sensory stimuli  $\mathbf{x}$  into intermediate embeddings  $\mathbf{z}$ . Subsequently, the “integrator” synthesizes these embeddings into the internal representation  $\mathbf{c}$ . Finally, the “predictor” network predicts future embeddings  $\mathbf{z}_{t+\Delta}$  from  $\mathbf{c}_t$  yielding the associated prediction error (PE). The learning objective is to minimize this PE. (b) Synthetic data generation paradigm with known latent variables. We generate motion trajectories by moving a 2D animal image, chosen randomly from one of eight categories (right), within a square arena (left). Each frame is characterized by the object position  $r_{1/2}$ , its scale to mimic position in the  $r_3$  direction, and the orientation angle  $\theta$ . For each trajectory we randomly select a constant velocity  $\mathbf{v}$  and angular velocity  $\omega = \dot{\theta}$ . To obtain pixel images that serve as network inputs, we further place these motion trajectories on a gray background and add Gaussian pixel noise, and elastic collisions at the borders (Methods). (c) Two example trajectories shown frame-by-frame. (d) Scatter plots of linearly decoded latent variables from learned internal representations  $\mathbf{c}$  for a representative sample of cat (blue) and dog (red) input images. Velocity ( $v_1, v_2$ ) and object identity  $S$  (left). Orientation ( $a, b$ ) and  $S$  (right). (e) Decoded velocity components  $\hat{v}_{1/2}$  as a function of ground truth on held-out data. (f) Same as (e) but for  $\hat{a}$  as function of  $\theta$ . The model (RPL) has formed untangled and structured representations of the latent variables. (g) Linear readout accuracy of animal categories  $S$  (left) and coefficient of determination ( $R^2$ ) (right) for position  $\mathbf{r}$ , velocity, orientation  $\theta = (a, b)$ , and angular speed  $|\omega|$  for linear decoders trained on the internal representations  $\mathbf{c}$  for  $n = 3$  independent networks, trained with the RPL and invariance learning (IL) objectives. Decoding performance for  $n = 1$  supervised end-to-end training (Sup.) and a randomly initialized network (Rand.) are given for reference. RPL decoding accuracy for object identity, and  $R^2$  values are close to the performance of networks trained end-to-end with supervised learning, whereas IL fails to accurately represent the velocity, orientation, and angular speed. (h, i, j) Same as d-f but for IL.

We next wondered whether an alternative asymmetric JEPA implementation, in which the prediction target is the integrator output rather than the encoder output, could also learn to extract both invariant and dynamic latent variables. If so, this would indicate that the asymmetric prediction character of a JEPA is key to learn both invariant and equivariant representations, while being insensitive to where the prediction errors are computed in the circuit, i.e., before or after the integrator. To investigate this question, we defined context-predictive learning (CL) as a third learning paradigm which minimizes prediction errors resulting from predicting the future representations at the level of the integrator  $c_t$  by minimizing prediction errors  $e_{t+\Delta}^{\text{CL}} = c_{t+\Delta} - f(c_t)$ . We then evaluated the representations using linear probing as before. Again, decoding accuracy from the resulting network representations was worse than RPL, but better than the IL objectives (Supplementary Fig. S2 [↗](#)). Hence, RPL's learning objective of predicting embeddings from a context-aware internal representation is crucial for learning untangled representations of invariant and dynamic latent variables.

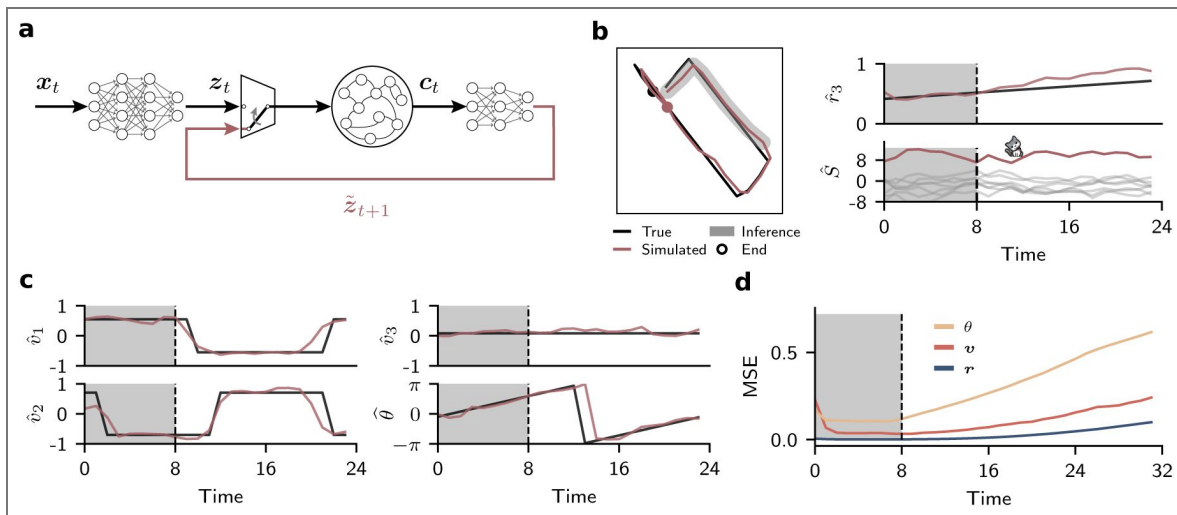
We speculated that CL's performance gap may be due to the integrator being able to internally “hallucinate” predictable signals that minimize the prediction error, even though they are unrelated to the input. We reasoned that this issue should become more pronounced when extracting latent variables from noisy input data, thereby requiring the network to rely more heavily on the integrator to infer them by integrating over time. To test this idea, we generated synthetic videos as before while systematically increasing pixel noise and frame dropout, i.e., introducing different numbers of consecutively blank frames at random locations in the video (Supplementary Fig. S6a [↗](#)). Consistent with the above hypothesis, we found that the performance gap between RPL and CL widened with increasing noise levels and number of dropped frames (Supplementary Fig. S6b,c [↗](#)), while RPL's overall performance remained relatively stable. Moreover, CL's performance variability between independently trained models was higher, suggesting that the CL paradigm is less robust. Thus, RPL's particular asymmetric architecture is essential to robustly untangle dynamic and invariant latent variables. Consequently, we will now focus on comparing RPL to IL for the remainder of the article.

## RPL circuits learn an internal world model

In addition to representing objects and their motion trajectories, animals need to anticipate future states of the world based on past observations. We wondered whether RPL implicitly learns a world model capable of such mental simulation. We reasoned that if the recurrent circuit indeed learns a world model, it should simulate motion trajectories using its internal representations in an autoregressive way. To test this idea, we introduced a dynamic switch in the network setup that routes the predicted embeddings  $\hat{z}$  back as input to the integrator, replacing the feedforward input from the environment (Fig. 3a [↗](#); Methods). Next, we primed a network in feedforward mode with the first eight frames of an animal motion stimulus to allow the model to infer the underlying latent variables. We then switched to feedback mode for 16 time steps. Throughout, we decoded all latent variables from the internal representations as before. We found that the networks in feedback mode simulated motion trajectories that closely resembled the position, scale, and object identity of the actual trajectory (Fig. 3b [↗](#); Supplementary Fig. S7 [↗](#)). Similarly, the velocity and orientation estimates resembled the ground-truth dynamics (Fig. 3c [↗](#)). Finally, to check how errors evolved over time, we calculated the average mean squared error over many simulated trajectories for different dynamic variables. This analysis showed that errors dropped rapidly during the inference phase, as expected, and slowly increased during the simulation phase (Fig. 3d [↗](#)). Together, these findings show that the learned internal representations can reliably simulate future object motion and that RPL learns a world model.

## RPL learns at different levels of abstraction

In the moving animal task studied above, the latent-variable dynamics were simple and predominantly linear, and each animal category was associated with a single image reducing the need for representation learning. We wondered whether RPL's specific asymmetric recurrent



**Figure 3. RPL circuits learn an internal world model.**

(a) Setup for simulating trajectories using internal representations. We added a dynamic switch allowing to replace feedforward input to the integrator with internal predictions (cf. Fig. 1). To simulate future trajectories, we provide the network with external-world sensory input during an initial inference period. Afterwards, we replaced feedforward input to the integrator with predicted embeddings in an autoregressive manner. (b) Example of a decoded trajectory (red) compared to the ground truth (black) in the  $r_1$ - $r_2$  plane (left) and  $r_3$  as a function of time (top right; cf. Fig. 2a) as well as the object decoder logits (bottom right). In this example we used an initial inference period of eight frames of a cat trajectory (shaded) followed by 16 steps of internal simulation. (c) Decoded velocity and orientation during the inference period (shaded) and simulation. The decoded values closely follow the ground truth. (d) Mean-squared error of the decoded quantities compared to ground-truth as a function of time. The error increases monotonically over simulation time.

architecture would enable it to simultaneously learn representations at different levels of abstraction for more diverse stimulus sets and more intricate sequence structures.

To investigate this question, we generated structured sequences of handwritten digits from the MNIST dataset [46]. Specifically, we generated the sequences using the following grammar. Sequences consisted of distinct digit triplets interspersed by zeros (Fig. 4a). Only triplets within the clusters of digits {1, 2, 3}, {4, 5, 6}, or {7, 8, 9} were allowed, covering all possible permutations within each cluster. To generate a sequence, we first selected a random cluster, sampled a random permutation within it, and then sampled random instances of the corresponding digits from the MNIST dataset (Fig. 4b). We continued by sampling the next cluster and triplets according to our defined grammar (Methods). We then fed the resulting image sequences to a network model for training with the RPL and IL learning objectives.

To probe the information encoded in the learned representations, we further trained three sets of linear classifiers to decode digit identity, cluster identity, and triplet identity (Fig. 4c; Methods). Encoding triplet identity requires a sense of temporal direction and a specific memory of recent digits. In contrast, cluster identity is direction-invariant and does not require distinguishing digits within a cluster. We thus expected RPL, an asymmetric JEPA, to yield high accuracy across all three tasks, in contrast to IL which, due to its symmetric objective, should only perform well at representing cluster identity. This is indeed what we found. RPL-trained network representations solved all three decoding tasks (Fig. 4d; Supplementary Table S2) with accuracy comparable to that of a fully supervised model.

In contrast, IL resulted in lower-fidelity representations of digit and triplet identity, on par with an untrained random network, while cluster identity was still encoded reliably as expected. Nonetheless, we also found that networks trained with IL could still represent individual digits somewhat faithfully at the encoder level by training a separate set of linear classifiers on the encoder outputs  $\mathbf{z}$  (Supplementary Note S1). However, this improved digit-decoding accuracy was still lower than for RPL. This analysis highlights again the importance of asymmetric prediction error computation through the recurrent pathway for learning time-resolved and abstract representations, even when these representations themselves do not depend on temporal context, as is the case for individual digit stimuli.

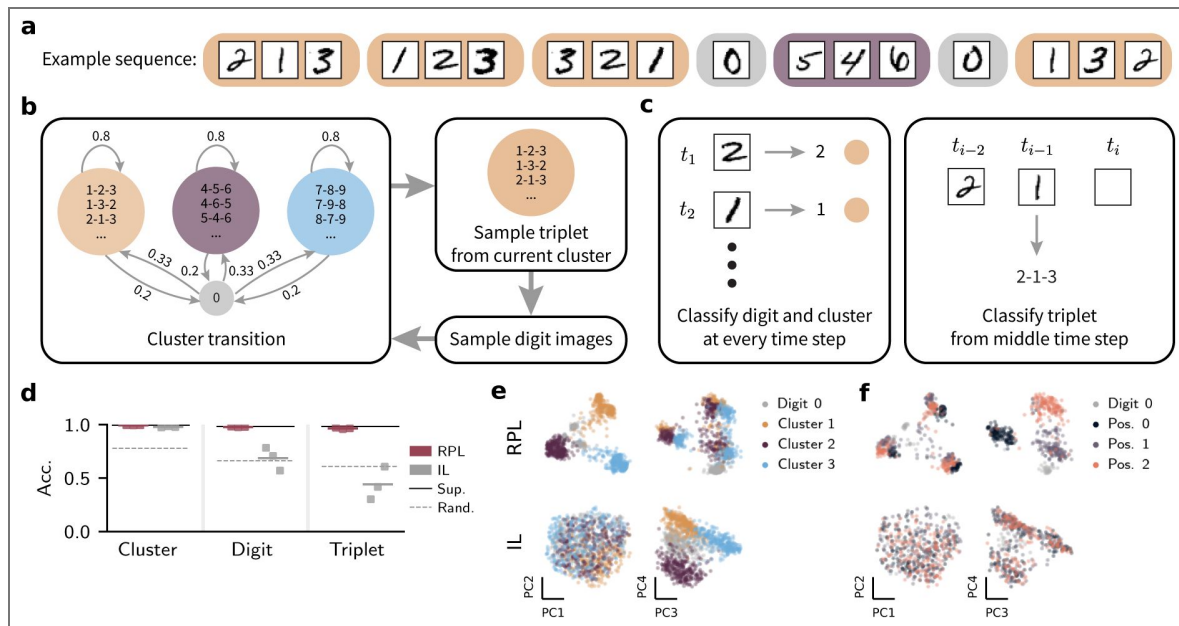
To gain further insight into how the temporal context is encoded in the trained networks, we visualized the principal components (PCs) of the learned representations (Fig. 4e,f; Supplementary Fig. S8). While cluster identity was clearly separated by both methods, e.g., along the first two PCs for RPL, and third and fourth PCs for IL (Fig. 4e), it was evident that RPL resulted in more finely structured representations. For instance, RPL representations encoded the temporal position within a triplet in the third and fourth PCs, essentially creating an internal clock that tracks sequence progression (Fig. 4f). Furthermore, PCs five through eight displayed a rich, though less interpretable, organization (Supplementary Fig. S8), presumably encoding the digit identity. In contrast, IL lacked such a fine-grained organization beyond clearly separating cluster identity, lending further support to our hypothesis that RPL's asymmetric architecture is essential for learning abstract sequence structure that requires a sense of direction, whereas IL's symmetric setup is limited to representing invariant features. Thus, RPL's asymmetric recurrent architecture is necessary to simultaneously learn complex representations at different abstraction levels from abstract sequence stimuli.

## Representation learning on real-world data

Having established RPL's effectiveness in extracting latent variables from synthetic datasets, we wondered to what extent these findings would generalize to real-world data. To address this question, we used published video data of behaving mice in an open-field arena [47] in which nose position  $\mathbf{r}_N$  and tail root position  $\mathbf{r}_T$  had previously been annotated with DeepLabCut [48]. Based on these annotations, we further estimated animal orientation  $\theta$  and velocity  $\mathbf{v}_T$  (Fig. 5a; Methods) as latent variables to evaluate the quality of the learned representations.

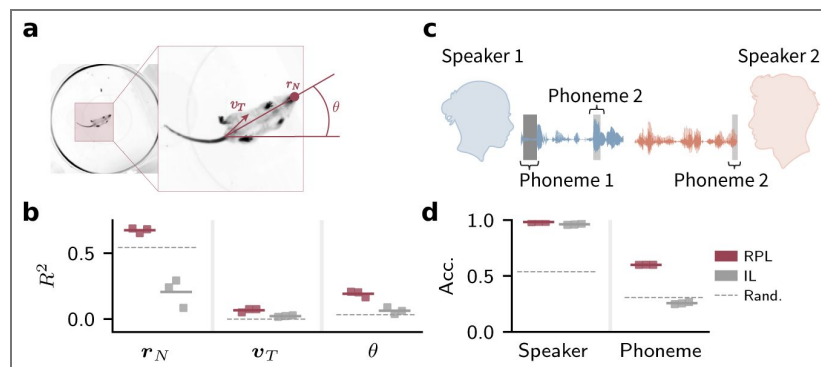
**Figure 4. RPL learns at different levels of abstraction.**

(a) An example sequence of digit triplets separated by zeros. (b) Schematic of the sequence generation paradigm. We generated triplet sequences containing randomly ordered handwritten digits sampled from one of three digit clusters. Transitions between clusters took place with probability 0.2 accompanied by a “zero”. For any given sequence, each digit instance was randomly chosen from the MNIST dataset. (c) Schematic of the tasks used to investigate the emergence of abstraction by asking how well the network representations encode cluster, digit, and triplet identity respectively. Solving these tasks requires both a mapping from images to abstract digit and cluster representations as well as knowledge about the transition structure. To quantify the representational quality, we trained a linear classifier for each task (Methods). (d) Linear decoding accuracy of the cluster, digit, and triplet identities for RPL and IL for  $n = 3$  independently trained networks. Decoding accuracy for end-to-end supervised learning (solid line) and a randomly initialized network (dashed line) are given for reference. RPL achieves high accuracy on all tasks, whereas IL results in a preference for cluster identity. (e) Principal component (PC) projections of the learned representations colored by cluster identity. Both RPL and IL clearly separate the clusters. (f) Same as (e) but colored by temporal position within a triplet. Only RPL yields representations of the abstract temporal transition structure.



**Figure 5. RPL extracts latent variables from real-world data.**

(a) Schematic of real-world video data of behaving mice in an open-field arena [47]. We used annotated key points to compute position  $r_N$ , velocity  $v_T$ , and orientation  $\theta$  as proxies of latent variables for evaluation (Methods). (b)  $R^2$  values of decoded position, velocity, and orientation for RPL and IL from  $n = 3$  independently trained networks. Decoding performance for a randomly initialized network (dashed line) shown for reference. (c) Schematic of real-world speech data from the Librispeech corpus [49]. We relied on speaker identity and phoneme labels to assess model performance. (d) Linear decoding accuracy of speaker identity and phoneme sequence for RPL and IL for  $n = 3$  independently trained models. Decoding from a randomly initialized network given for reference (dashed line). RPL has high decoding accuracy for speaker identity and phonemes. In contrast, IL yields high accuracy only for the speaker identity.



We then trained RPL networks on the videos and used linear probing for evaluation as before. For comparison, we also trained the same network using the IL objective. We found that decoding from RPL representations yielded higher  $R^2$  values across all variables (Fig. 5b), consistent with previous findings, that RPL untangles representations during learning from real-world videos. We observed that velocity was more difficult to decode, possibly because mice remained stationary in large portions of the videos, leading to a velocity distribution with a peak near zero (Supplementary Fig. S9).

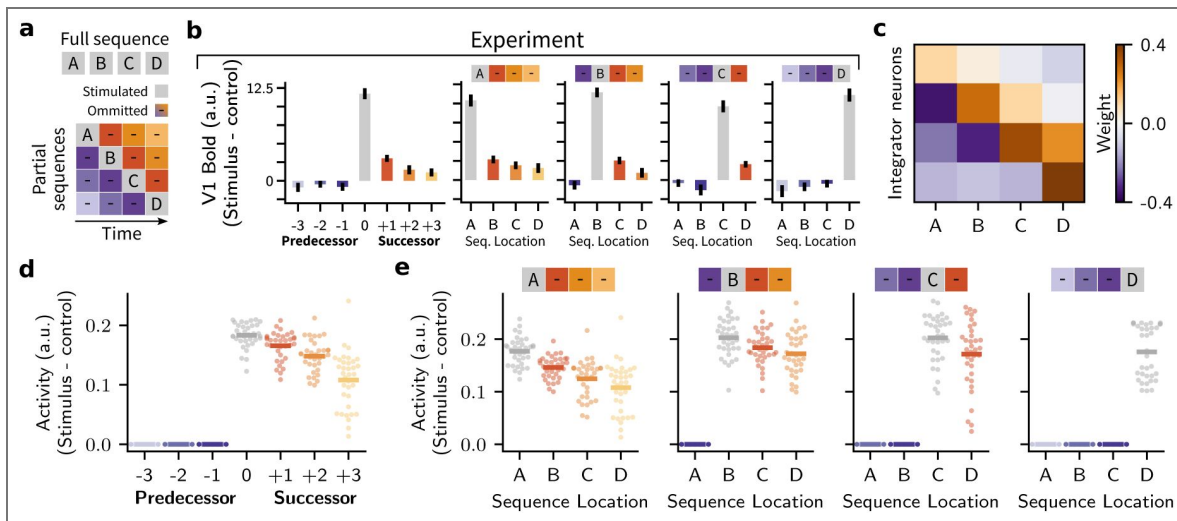
Since the mouse videos were structurally quite similar to our synthetic data above (cf. Fig. 2), we wanted to see whether RPL would yield similarly useful latent representations in a different data domain. To investigate this question, we trained both RPL and IL models on spoken-language recordings from the Librispeech corpus [49]. This dataset includes labels for speaker identity and the phonemes of the spoken text (Fig. 5c). Since phonemes change within and across sentences whilst speaker identity remains unchanged, we expected that RPL would learn representations that allow decoding both, whereas symmetric IL objectives would focus on speaker identity. To check whether this is the case, we evaluated the learned representations by linearly decoding speaker identity and phoneme sequence (Methods). Our results confirmed our expectation. RPL yielded higher accuracy for both speaker identity and phonemes, whereas IL failed to represent phoneme information (Fig. 5d). Thus, RPL extracts meaningful latent variables across different real-world data domains.

## RPL learns successor-like representations comparable to V1 in humans

We next wondered to what extent RPL network models predicted *in vivo* data. To that end, we first focused on signatures of successor representations, whereby a neuronal population encodes a predictive representation of future states, such that each state representation reflects the expected future occupancy of other states under a given policy [50]. Recent work described successor-like representations in fMRI recordings from the primary visual cortex of humans [45]. During these experiments, humans were exposed to a sequence of four non-overlapping items (A-B-C-D) across multiple trials. After an initial learning phase, the participants were exposed to either the entire sequence or only one of the four sequence items, referred to as partial sequences (Fig. 6a). The data from V1 indicated activation at the successor but not the predecessor locations upon partial sequence exposures (Fig. 6b).

We wanted to know whether RPL networks exhibit similar activation patterns at successor locations after being exposed to comparable sequential stimuli. To investigate this question, we trained minimal RPL circuit models ( $n = 35$ ; Methods) on an (A-B-C-D) sequence to mirror the experimental protocol. After sorting the neurons in the integrator based on their selectivity, we visualized the average weight matrix associated with the neural selectivity to each sequence item (Fig. 6c). This matrix showed the largest entries on the main diagonal, positive values on the two upper off-diagonals, and negative values below the main diagonal, a clear signature of successor representations.

We then computed the average neural activity at the stimulated, predecessor, and successor locations in response to partial sequences, after removing baseline neural activity from all neural responses. We observed that, as in the fMRI recordings [45], the stimulated location showed the highest activity, while activity at successor locations followed a decaying pattern. In contrast, there was no activity at the predecessor locations consistent with the experimental data (Fig. 6b,d). Similarly, we plotted neural activity in response to each partial sequence to see whether this pattern also existed in selective neurons for non-stimulated sequence items at successor locations. We observed a consistent pattern of activity at stimulus and successor but not predecessor locations (Fig. 6e). Thus, RPL entails successor-like representations as observed in V1 of humans experiencing repeating sequence stimuli.



**Figure 6. RPL learns successor-like representations comparable to human V1.**

**(a)** Schematic of the full and partial sequences used in the experiment. **(b)** Relative BOLD signal responses reproduced in human V1 [45]. Relative stimulus response to partial sequence stimuli (left). Individual responses to partial sequences (right). Partial sequence stimulation caused brain activity encoding successor locations but not predecessor locations. Error bars represent the standard error of the mean. Data extracted and reproduced from Ekman et al. [45]. **(c)** Relative neuronal selectivity matrix extracted from the RPL model representations ( $n = 35$  models). The negative values below the diagonal indicate inhibition of activity at predecessor locations. Positive values above the diagonal indicate that the successors' representations are driven by the stimulus. **(d)** Average neural activity for all partial sequences at non-stimulated and stimulated locations after removing the baseline activity as in (b). Each dot is a trained model and the solid lines represent the average activity over all models. The plot indicates positive neural activity at stimulated and successor locations, and no activity at predecessor locations. **(e)** Neural activity at all sequence locations in response to each partial sequence. These plots indicate the decaying positive response at successor locations for individual partial sequences, hinting at successor-like representations in response to each partial sequence.

## RPL learns abstract sequence representations similar to macaque PFC

Having observed that RPL learns predictive representations consistent with those in a primary sensory cortex, we asked whether our model could also capture more abstract representations ascribed to higher-order cortical areas, such as the PFC. A widely used experimental setup to investigate abstract sequence learning is the local-global oddball paradigm (Fig. 7a). Here, an animal is exposed to trials or sequence presentations consisting of four stimuli one after the other, either repeating the same stimulus (xxxx or xx for short) or an oddball in the last slot (xxxY or xY for short), constituting a local deviant. Repeating trials are organized into blocks, in which the animal is conditioned to expect a particular sequence type, e.g., xx or xY, as the standard by dint of frequent presentation at the beginning of the block. Later trials in the block will be either global standards or infrequent global deviants, depending on whether they match the context set at the start of the block. For example, in a block where xY is the standard, an xx trial denoted (xx|xY) constitutes a global deviant but not a local one.

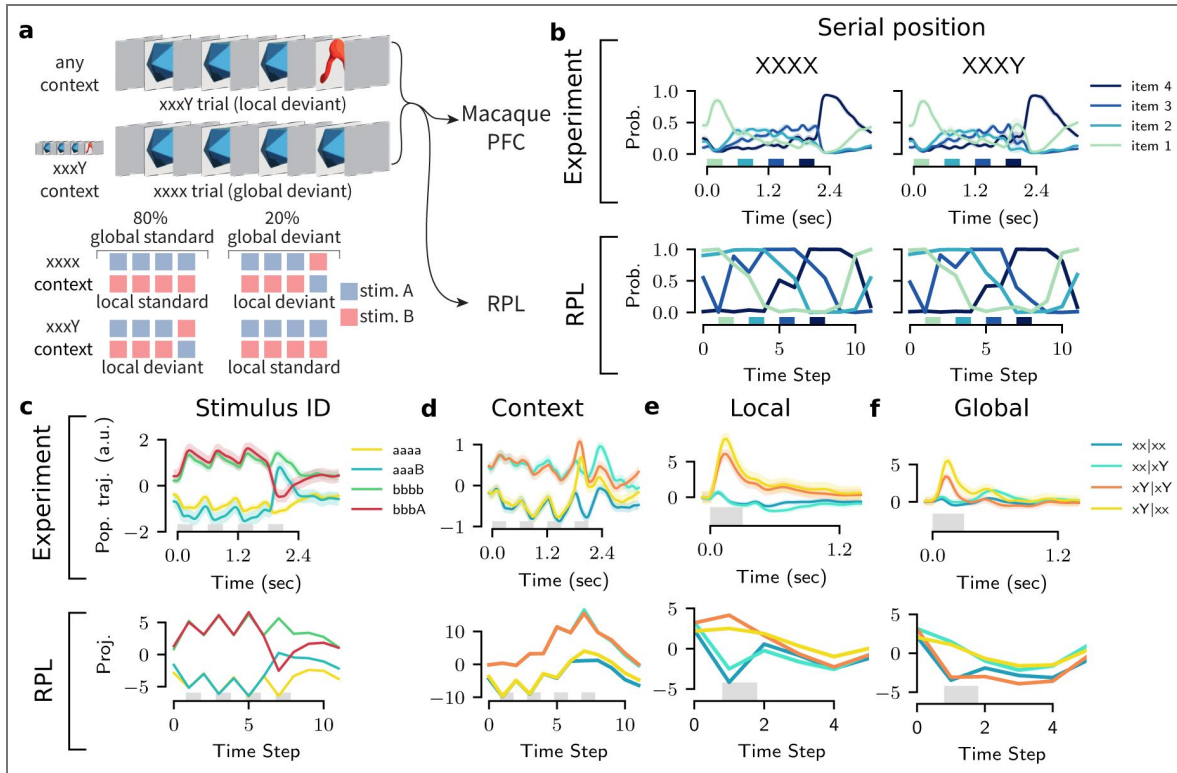
One such experiment was performed by Bellet et al. [44] who recorded neurons from the macaque PFC while exposing them to visual local-global oddball sequences. They found that populations of neurons in the PFC spontaneously learned to encode several abstract sequence properties, such as their serial position within a trial of four stimuli, stimulus identity, global context, and the presence of local and global deviants.

Since we had already observed that RPL is capable of abstract sequence learning (cf. Fig. 4), we hypothesized that RPL should also be able to learn the sequence properties of local-global oddball sequences. To test this hypothesis, we trained an RPL circuit model on the identical visual stimulus sequences that Bellet et al. [44] presented to the macaques, closely mirroring their experimental protocol. Then, following their analysis (Methods), we extracted neuronal population vectors by training linear readouts to decode: serial position, stimulus identity, global context, and local/global deviance.

We found that the serial position of an item in a sequence and stimulus identity could be decoded reliably (Fig. 7b,c) as in the experiments, albeit without the experimentally observed accentuation of the first and the last item, which are presumably due to primacy and recency effects in the experiment [44]. To test whether the RPL network had learned an internal model of the sequence structure, we decoded the global context variable from the internal representations. As observed in the authors' original analysis, we identified neural subspaces in our model that coded for the global context inferred from previous sequence observations (Fig. 7d). Finally, we assessed our model's responses to local and global deviants. We found qualitatively similar separation of population responses corresponding to the presence of local deviance in both predicted (xY|xY) as well as unpredicted (xY|xx) local deviant stimuli (Fig. 7e). Similarly, the model showed an explicit encoding of the presence of global deviance (xx|xY and xY|xx; Fig. 7f) in its population responses, as in the experiment. Although our model did not exhibit the onset and primacy effects observed in the data, presumably because it lacked neuronal adaptation mechanisms, we found that the representations learned via RPL qualitatively matched the relative encoding in PFC. Thus, RPL learns abstract sequence representations whose relative encoding structure is consistent with PFC activity.

### Untangling in a cortical model with local prediction error circuits

So far, we have used deep encoder networks, computing prediction errors at the top-level output, using backpropagation of error (backprop) to update all weights. However, experimental evidence points to local mismatch computation across the hierarchy as early as primary sensory cortex [12]. Moreover, how the brain approximates backprop to solve the credit assignment problem remains an open question [51]. Hence, we wondered whether RPL remains effective without end-to-end optimization and to what extent.



**Figure 7. RPL learns abstract sequence representations similar to macaque PFC.**

**(a)** Schematic of the local-global oddball paradigm used by Bellet et al. [44]. Example trials consisting of four sequential stimuli show an “xxxY” and an “xxxx” trial in the context of frequent “xxxY” presentations (top). A given trial can be a local standard or deviant depending on the local sequence structure of the four stimuli and a global standard or deviant depending on the context set for a block of several trials (bottom). Figure adapted from [44]. **(b)** Comparison of temporal serial position representation in experiment (top) and the model (bottom). Plots show best linearly decoded probabilities of the four serial positions (items 1 through 4) averaged over all xxx and xxxY trials. RPL learns representations that encode serial position more robustly than the macaque PFC. The experimental data corresponds to Monkey A and has been reproduced from publicly available data and code. **(c)** Comparison of stimulus identity representation in the experiment (top) and the model (bottom). Plots show projections of neural activity or RPL representations on directions that best separate pairs of input stimuli (Item a vs b) that comprise x and Y, averaged across trials (Methods). RPL representations encode stimulus identity similar to the experiment. **(d)** Same as (c) but for representation of the global context set by the frequent trial. **(e)** Same as (c) but for the presence or absence of a local deviant at the fourth sequence position. **(f)** Same as (c) but for the presence or absence of a global deviant, decoded from the fourth sequence position.

© 2024, Bellet ME. Panel A was reprinted with permission from Figure 1 Bellet et al. 2024 [44], which was published under a CC BY-NC-ND <https://creativecommons.org/licenses/by-nc-nd/4.0>. Further reproductions must adhere to the terms of this license.

To investigate this question, we built a hierarchical RPL (h-RPL) model in which we implemented different cortical areas as local RPL circuits that learn independently. Functionally, each circuit consists of a shallow encoder, an integrator, and a predictor putatively implemented by distinct cell types. One possible mapping of the model onto cortical anatomy is to view the encoder as being implemented through lateral connectivity between adjacent areas, whereas vertical inter-laminar connectivity implements the integrator and predictor, resulting in a circuit that produces local prediction errors (fig. 8a) reminiscent of mismatch responses observed in superficial cortical layers [12].

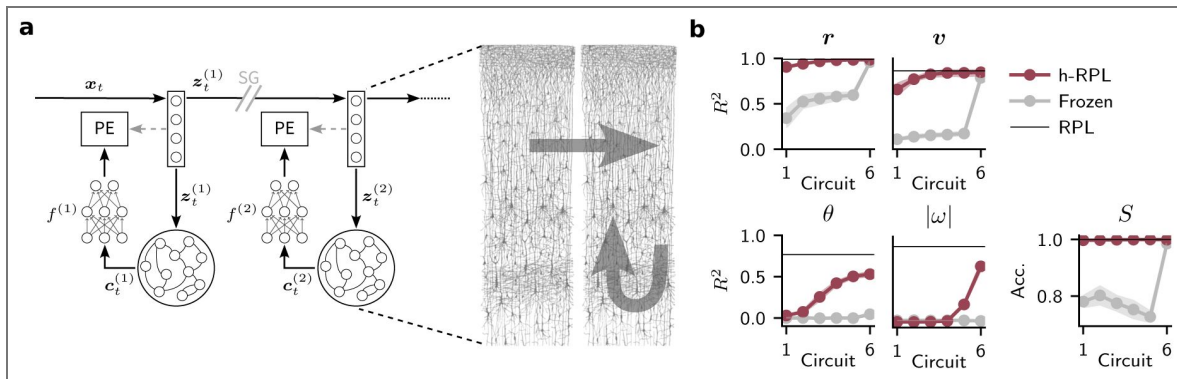
We trained this h-RPL model on synthetic moving animal data (cf. Fig. 2). For simplicity, we trained every circuit using backprop, but we prevented gradients from flowing between circuits across the hierarchy (Methods). Doing so confines credit assignment to individual circuits with a few synaptic transmission steps, a setting in which biologically plausible backprop approximations are effective [53–56]. In consequence, each circuit  $l$  learned independently in this hierarchical network model.

To assess the quality of the learned representations, we trained linear decoders for each circuit and extracted the latent variables as before (Fig. 8b). This analysis revealed that the linear decoding accuracies and  $R^2$  values for the dynamic variables increased with depth along the hierarchy, plateauing at performance levels comparable to those of the end-to-end trained RPL setup. To check whether this linear decodability of the evolving latent variables also implied “perceptual straightening,” a phenomenon characterized by the progressive transformation of curved stimulus trajectories into more linearly evolving representations along the visual cortical hierarchy [57, 58], we further quantified straightness along the model’s hierarchy. We found, that straightness indeed increased along the hierarchy and was due to learning (Supplementary Fig. S10).

One potential caveat with this analysis is that the hierarchical encoders jointly constitute a deep convolutional network with increasing receptive field sizes. To check that the observed effects were not merely due to the progressively larger receptive fields, we trained another instance of the same hierarchical model in which we kept the weights in all but the last circuit of the hierarchy frozen. This resulted in an overall reduction in linear decodability across the hierarchy (Fig. 8b; Supplementary Figure. S10). While the drop in  $R^2$  for position, velocity, and object category could be almost entirely compensated for by the optimized last circuit, this was not the case for orientation and angular speed. Additionally, we observed a steady increase in the autocorrelation time of the representations over the hierarchy for h-RPL but not in the frozen counterpart (Supplementary Fig. S11). These findings are consistent with the idea that better decodability across circuits is not merely a result of increasing receptive field size, but rather the progressive untangling of representations across the hierarchy. Thus, RPL is a powerful learning algorithm for finding dynamic latent variables even in a hierarchical setting, which resonates with the core ideas of predictive processing and the experimentally observed mismatch neurons in primary sensory cortices.

## Discussion

We introduced recurrent predictive learning (RPL), a self-supervised learning model that learns untangled representations of objects and their dynamics directly from sequential observations alone, without requiring labels, an implausible decoder, or a reconstruction objective. Formally, RPL constitutes a JEPA that extracts both invariant and equivariant representations across multiple levels of abstraction from synthetic and real-world data. Beyond representation learning, RPL also acquires an internal world model that can simulate plausible future states purely based on its internal representations. Crucially, the learned internal representations resemble neuronal activity patterns reported in macaque PFC and exhibit successor-representation-like signatures comparable to those observed in human primary visual cortex.



**Figure 8. Untangling in a cortical model with local prediction error circuits.**

(a) Schematic of the network model consisting of six identical RPL circuits in a hierarchy, each containing a shallow encoder, an integrator and a predictor network. Each circuit computes its own prediction error while Stop Grad (SG) operations prevented gradients from flowing between the individual circuits. In this h-RPL model horizontal intra-laminar connectivity implements the encoder network. In contrast, vertical connectivity accounts for the integrator and predictor circuitry. Cortical anatomy drawing (right) adapted from Cajal [52]. (b) Decoding  $R^2$  values and accuracy of the latent variables (cf. Fig. 2b,c) for each RPL circuit in h-RPL (red) and a variant in which only the last circuit in the hierarchy was trained (gray). Shading corresponds to one standard deviation over  $n = 3$  independent models. The solid black line denotes the performance of RPL with a deep encoder (cf. Fig. 2j). The h-RPL model progressively learns latent variable representations comparable to end-to-end trained RPL.

Learning in RPL is driven by internal prediction errors generated within a specific circuit motif: an encoder with little working memory feeds into an integrator-predictor circuit that integrates temporal information and learns to anticipate future encoder embeddings. This motif is functionally universal, enabling hierarchical networks with local predictor circuits to progressively uncover latent variables. While invariant representations can be learned by local learning rules in conjunction with appropriate decorrelation mechanisms [11], we showed in this work that learning equivariant representations requires a circuit with distinct functional elements for encoding, integration, and prediction. We expect analogous functional motifs to exist in brain areas that rely on self-supervised learning, implemented through dedicated cell types and connectivity patterns that jointly shape plasticity. In this way, our work provides a conceptual blueprint for interpreting neuronal circuits, while making concrete, testable predictions.

## Relation to previous modeling work

The notion that the brain builds and progressively updates an internal model of the world through comparison with observations is a central tenet of systems neuroscience [10, 12, 13, 17, 19, 27, 28, 59–64]. Classical predictive coding models based on this idea assume an underlying generative model that *reconstructs* sensory input [13, 15, 16, 18, 65]. The generative model is learned and updated by computing prediction errors in the input space. However, classical predictive coding is under increasing pressure due to its difficulty in explaining recent experimental data [25, 26, 28]; moreover, it yields poor representations for discriminating between invariant concepts. Modern variants of predictive coding alleviate this problem by combining a supervised loss function with the classic unsupervised reconstruction loss [24]. Still, doing so requires labels. These difficulties might stem from its core assumption of an underlying generative model. Indeed, recent studies argue that reconstruction, a core component of generative models, is inefficient for learning high-level abstract representations because it tends to focus on behaviorally irrelevant low-level stimulus details [42, 66].

In contrast, JEPAs like RPL depart altogether from the idea of an underlying generative model. Nevertheless, RPL, and h-RPL in particular, share some essential elements with classical hierarchical predictive coding models. In all cases, learning is driven by prediction errors that are computed locally and used to update model parameters and internal representations, a mechanism consistent with the existence of prediction-error-like neurons in sensory cortex [12, 27, 67]. In this respect, h-RPL aligns with established predictive coding models in which cortical circuits continuously generate and refine predictions about incoming signals [13, 15, 68].

However, a central distinction lies in what is communicated between the circuit elements. Classical predictive coding architectures rely on the explicit propagation of prediction errors across hierarchical levels, typically assuming a functional segregation between representation neurons and error neurons. In contrast, the JEPA-like organization underlying h-RPL does not require prediction errors to be transmitted beyond the local circuit in which they are computed.

Instead, separate network structures embed related signals and learn to predict each other's latent states via locally computed errors. As a result, inter-areal communication is dominated by representations rather than error signals, with learning occurring through the mutual alignment of internal embeddings. This shift fundamentally changes the representational objective of the network, emphasizing abstract, predictive latent variables over the reconstruction of sensory input. Finally, we should note that this shift is not mutually exclusive with modern forms of *discriminative* predictive coding, which, in conjunction with a predictive self-supervised loss as discussed below, could support representation learning by approximating backprop [24, 55].

Based on these arguments, this article advocates for a revised predictive processing framework, one that avoids explicit generative models, meaning that it neither relies on reconstruction in the input space nor supervised objectives, and instead relies on SSL via prediction in the representation space. Recent studies have taken steps in this direction by developing reconstruction-free predictive SSL models that compute predictions directly in the embedding or representation space. For example, Ororbia et al. [69] proposed a model that relies on spatial prediction across visual scene embeddings. Other models rely on contrastive learning [10, 70] or

temporal prediction [11] across sequences of augmented images reminiscent of slow feature analysis [71]. These studies have provided valuable insights into learning invariant representations without supervision. However, they primarily focused on category-level invariance induced by hand-crafted augmentations and did not explicitly address how biological agents may simultaneously acquire equivariant representations that preserve behaviorally relevant information, such as motion direction or relative spatial relationships.

Related to our work, Nejad et al. [20] introduced a temporal predictive SSL model with a putative mapping onto neocortical layers. In this model, prediction occurs in the latent space spanned by Layer 5 (L5). Specifically, Layers 2/3 (L2/3) predict the next L5 representation. To avoid representational collapse, L5 is either trained like an autoencoder with an additional reconstruction objective, or by adding a variance-maximization criterion. The resulting model learns to predict orientation changes while capturing several experimental observations. In contrast to RPL, the study used shallow encoder networks and utilized sequential Gabor patches as sensory input, with the rotation direction provided as top-down input. Similarly, Asabuki et al. [72] investigated predictive sequence learning in a local recurrent circuit model, emphasizing biologically plausible learning rules. While these studies make notable contributions toward understanding how brains may implement predictive learning, they focused on comparatively simple network architectures and elementary learning tasks. They did not address the emergence of rich, abstract representations from complex sensory streams.

While the above studies similarly argue for a set of core principles that should underpin a biologically grounded theory of representation learning, i.e., one that does neither rely on reconstruction nor supervision, the same principles naturally align with the class of SSL models developed in machine learning under the umbrella of JEPAs [34, 36, 73, 74]. JEPAs operationalize the idea that internal representations can be learned by making different views or temporal segments of sensory experience predict one another, thereby preserving the core intuition of predictive processing while avoiding many of the constraints imposed by explicit generative models. Crucially, RPL is a JEPAs, which means it neither relies on a reconstruction objective nor on negative samples, and by avoiding the issues associated with generative models, it learns useful abstract representations from a stream of sensory experience, all the while remaining consistent with the core idea of predictive processing and neural circuit principles.

Although a plethora of different JEPAs have demonstrated their effectiveness in machine learning [11, 33, 36–42, 66], not all of them are biologically plausible. For instance, existing studies made strong assumptions about the allowed transformations between consecutive stimuli [39] or were not applied to real-world problems and lacked a circuit-level interpretation [75]. In part, this lack stems from current JEPAs typically relying on biologically implausible transformer networks, predicting randomly masked parts of the stimulus, and optimizing a single global prediction error during training, all of which are difficult to reconcile with neural anatomy. Finally, the majority of research has focused on invariance learning, in which representational quality is assessed primarily through object classification tasks, e.g., via linear decoding or probing. Such quantification of model performance at the level of internal representations is a welcome and necessary departure from evaluating representational quality solely on a model's ability to reconstruct stimuli at the pixel level, as commonly done for classical predictive coding models [15, 76]. However, focusing solely on invariance learning seems like an oversimplification. An increasing number of machine learning studies are acknowledging the importance of equivariant representations [33, 39, 77–80], albeit with the limitations mentioned above and the assumption of direct access to the underlying transformations.

RPL avoids several implausible aspects of the above studies. On the one hand, RPL is a temporal JEPAs based on recurrent neural networks rather than transformers. On the other hand, it does not rely on masking as commonly used in self-supervised video models [40, 41]. Velarde et al. [81] studied a related recurrent JEPAs while focusing on biologically motivated learning rules that avoid backprop through time. While similar to our approach, their setup closely resembles our CL objective, which proved less effective in our study. Finally, h-RPL is, to the best of our knowledge, the first implementation of a hierarchical JEPAs that optimizes multi-level prediction errors as

initially proposed by LeCun [35]. Functionally, RPL learns both invariant and equivariant representations consistent with the idea of active filtering [82], as well as a latent generative model that enables it to simulate plausible dynamics in the latent space when queried with its own predictions. Interestingly, the necessity of separating the feedforward prediction target from the integrator-predictor branch directly stems from this reliance on a plausible recurrent neural network architecture, which opens the door to further investigation of the putative links between our model and the cortical microcircuit.

## Relation to experimental observations and circuit structure

Our work reveals several points of convergence between RPL and experimental observations of cortical computation, spanning both early sensory and higher-order areas. At a functional level, RPL develops representations that are explicitly optimized for prediction over time, rather than reconstructing sensory input. In particular, RPL learns successor-like sequence representations, in which current stimuli evoke activity at future but not preceding locations, closely paralleling successor representations recently reported in human primary visual cortex [45] and broadly consistent with sequential activity recall observed in mouse V1 [83].

Beyond early sensory areas, RPL also acquires abstract sequence representations that resemble those observed in macaque PFC during local–global oddball paradigms [44]. These representations encode relational and temporal structure rather than sensory detail, consistent with the idea that higher cortical areas operate on increasingly abstract latent variables. More broadly, our results align with experimental and theoretical work suggesting that cortical representations are shaped by predictive objectives [14, 22, 84] and become progressively more invariant over the cortical hierarchy [85].

In h-RPL, we observed a similar systematic increase in temporal stability across circuits (Supplementary Fig. S11 [\[link\]](#)), mirroring timescale hierarchies observed across cortical areas [86]. Finally, h-RPL offers a mechanistic account of perceptual straightening phenomena in visual cortex [57, 58], which refers to the progressive reorganization of sensory representations such that temporally structured or correlated stimuli come to be represented as more stable, linearized trajectories in neural state space, improving predictability and downstream readout. Our work suggests that such effects may arise naturally from representation-space predictive learning rather than from explicit sensory reconstruction.

In relation to cortical anatomy, the central contribution of our model is not a one-to-one mapping between algorithmic components and specific cortical cell types or layers, but the proposal that representation learning of invariant and equivariant representations requires a circuit with specific elements: encoder, temporal integrator, predictor, and prediction-error computation, all implemented by different circuit elements and, ultimately, different neuronal populations. This separation of function is a defining feature of RPL and of JEPAs more generally. Importantly, it does not require that these components be realized in a single, unique, or canonical circuit layout.

In the minimal model studied here, these elements are arranged in a feedforward–recurrent motif without explicit top-down connections or recurrent interactions between circuit modules. This abstraction is an oversimplification of cortical reality, which was chosen deliberately to isolate the core circuit elements for predictive learning using internal representations, not to provide a complete account of cortical anatomy. The cortex contains dense recurrent connectivity within and across layers, extensive lateral interactions, and bidirectional long-range projections, any of which could support alternative implementations of the same underlying computation. As a result, the precise correspondence between model components and cortical microcircuits remains inherently ambiguous.

That said, an appealing idea that is broadly consistent with our h-RPL model is that the encoder could be implemented through lateral intra-layer connectivity, whereas the circuits required for generating predictions and prediction errors follow a vertical columnar organization. In this context, L2/3 neurons, which have been associated with prediction-error-like signals, are a prime candidate for computing prediction errors. Still, how the integrator and predictor map onto the

cortical circuitry is less clear. Even if these elements did follow a vertical organization, they could well be intermingled without clear laminar division. Finally, while our RPL implementation assumes a single encoder, the JEPA framework naturally admits alternative configurations with distinct encoder networks [36], which could reside in the superficial and deep layers of the cortex. This ambiguity underscores that the framework specifies which computations must be present, rather than where they must reside, and highlights the need for future experimental work to disambiguate among these possible biological realizations.

However, once a suitable set of candidate circuit elements for encoder, integrator, and predictor are identified, the model leads to a concrete prediction. In particular, acute perturbations that disrupt the coupling between encoder and integrator components should primarily alter the representations formed by the integrator–predictor circuit, while leaving the activity patterns associated with the encoder largely preserved. By contrast, chronic disruptions of this coupling are expected to interfere with representation learning itself, preventing the gradual alignment of internal representations and the emergence of stable, abstract latent variables.

Beyond the neocortex, the proposed separation between encoding, integration, and prediction naturally invites comparison with hippocampal function. A substantial body of work has characterized the hippocampus as a predictive system that learns temporally extended representations of experience, often formalized in terms of successor representations or predictive maps of future states [87–89]. In this view, hippocampal activity reflects expectations about future sensory or cognitive states rather than a reconstruction of current input, consistent with the representational emphasis of predictive learning. Recurrent dynamics within hippocampal circuits, particularly in CA3, together with sequential activity observed during navigation and replay, have been interpreted as supporting temporal integration and prediction [19, 90]. From this perspective, hippocampus and cortex may implement complementary predictive architectures operating at different temporal and spatial scales, with hippocampal representations providing rapidly learned, temporally extended predictions that could scaffold slower cortical representation learning [91, 92]. While we did not explicitly model hippocampal circuitry, RPL's emphasis on internal prediction aligns naturally with theories that treat the hippocampus as a generator of predictive representations rather than a passive memory store [93].

Overall, our results suggest that experimental observations traditionally interpreted through the lens of hierarchical predictive coding may instead reflect a more distributed predictive architecture, in which different circuit elements specialize in complementary roles. JEPAs provide a useful conceptual scaffold for organizing these roles, while remaining compatible with multiple circuit layouts and known deviations from strictly hierarchical organization. The model presented in this article gives one possible instance, but it provides also a blueprint outlining which minimal computations must be present. Where these reside, and how precisely they are implemented remains an open question. Elucidating the architectural motifs the brain employs, and under what conditions, will require targeted experimental and theoretical work.

## Limitations and future research directions

Our work has several limitations that we plan to address in the future. First, the link between RPL and concrete circuits in the brain is under-constrained. One key simplification of our model is the lack of top-down connections. Concretely, we assumed that the predictor in the h-RPL model is entirely local, with no top-down input, a crude simplification of cortical anatomy where top-down connections are pervasive and considered essential to provide lower sensory areas with high-level abstract representations [94]. Adding such connections would allow modulating predictions within each circuit with global context from higher-level circuits broadly consistent with experimental findings on the role of feedback in sensory cortices [25, 28, 64]. A similar combination of a feedforward sensory pathway with a separate recurrent feedback pathway has been previously proposed as the “dual counterstream” architecture, an organizational principle of cortical structure [95]. Exploring this connection and its functional implications is an exciting avenue for future research.

Another limitation of our model is that we relied on backprop, a biologically implausible training method. While we have shown that circuit-local training can be effective with h-RPL (Fig. 8), we still used backprop to train each circuit. In the future, we will investigate how biologically plausible credit assignment strategies could make the learning algorithm more credible [51, 55, 56, 96–101].

Moreover, latent-variable learning in RPL was sensitive to the prediction time horizon and the underlying distribution of the latent variable. In this article, we have only explored this connection to a limited extent (Supplementary Fig. S9). A systematic investigation of how the time horizon and latent distribution influence the learned representations will be the goal of future research.

Finally, we relied on linear decodability as a proxy to quantify the quality of learned representations. To that end, we extended the commonly used linear classification strategy to linear regression of dynamic latent variables. While we saw that RPL embeds unrelated latent variables in orthogonal subspaces and improves overall linear decodability, a key question is why. In fact, there is no immediately apparent reason why the learned representations should be embedded on linear manifolds, and presumably they are not. One reason behind the improved linear decodability may be that linear decoding benefits from the orthogonal embedding of unrelated factors, combined with the high embedding dimension. An important question for future research is thus to gain a deeper theoretical understanding of how representations of continuous latent variables emerge in self-supervised learning JEPAs and which inductive biases lead to representations conducive for their linear decodability.

Extracting abstract representations of the world's invariant and dynamic features is a crucial function of the brain that enables intelligent behavior. Here, we have presented a plausible computational model of self-supervised representation learning for both object identity and dynamic latent variables. Our framework extends a long tradition of predictive processing models while going beyond classical predictive coding and pixel space prediction. Specifically, our work builds a conceptual bridge between neuronal circuits and JEPAs, a promising machine learning architecture that does not rely on a reconstruction objective. While many important questions remain open, we believe that our model provides a fresh approach to understanding how the brain forms abstract internal representations from continuous sensory experience, and that exploring the research avenues outlined above holds great promise for gaining deeper insights into the computational principles underlying biological intelligence.

## Methods

In this article, we studied different representation learning paradigms that rely on prediction errors computed in internal representation space. We quantified their ability to extract latent variables like object identity, position, velocity, and orientation from streaming data such as video or audio. In all cases, we used deep neural networks comprised of three distinctly connected network elements: a feedforward encoder network, a recurrent integrator network, and a feedforward predictor network.

### Learning objectives

All networks were trained using backprop to optimize one of the following learning objectives defined on the internal representations.

### Recurrent predictive learning (RPL)

While most predictive processing models optimize for reconstructing stimuli in input space, RPL optimizes a predictive objective function defined entirely on internal representations. Concretely, RPL is a particular JEPA which predicts the embeddings  $\mathbf{z}_{t+\Delta}$  of future stimuli using a predictor network  $f$  that acts on the representations  $\mathbf{c}_t$  from the integrator (cf. Fig. 2a). So, the RPL objective is defined as:

$$\mathcal{L}_{\text{RPL}} = \|f(c_t) - \text{SG}(z_{t+\Delta})\|^2, \quad (2)$$

where SG corresponds to the Stop-Gradient operator. The RPL objective avoids representational collapse by combining the SG with an asymmetric prediction architecture and, consequently, does not require explicit variance regularization [37, 102, 103].

## Invariance learning (IL)

Invariance learning is an alternative approach for representation-space predictive learning that uses a symmetric prediction scheme with explicit regularization terms [11, 32] for preventing collapse. This is also sometimes referred to as a joint-embedding architecture [35, 38]. The IL objective defined on the output of the predictor network in an architecture identical to the one used in RPL (cf. Eq. (2)) is one such approach. We use this, along with a regularization objective  $L_{\text{reg}}$  similar to the one used by Halvagal et al. [11], as a baseline for comparison with RPL. The overall IL objective function is then defined as:

$$\begin{aligned} \mathcal{L}_{\text{IL}} &= \|f(c_t) - \text{SG}(f(c_{t+\Delta}))\|^2 + \mathcal{L}_{\text{reg}}(f(c_t)) \\ &= \|f(c_t) - \text{SG}(f(c_{t+\Delta}))\|^2 - \lambda_1 \frac{1}{M} \sum_{i=1}^M \log \sigma_i^2 + \lambda_2 \sum_{i=1}^M \sum_{j \neq i}^M C_{ij}^2, \end{aligned} \quad (3)$$

where  $\sigma^2$  is the overall variance of the activity of neuron  $i$  in the output layer (i.e.  $f(c_t)$ ),  $C_{ij}$  is the correlation coefficient between the activities of neurons  $i$  and  $j$ , and  $\lambda_1$  and  $\lambda_2$  are constants that control the strength of each regularization term. The two terms in  $L_{\text{reg}}$  prevent the neurons from activity collapse and learning identical information, respectively.

Other variations of the IL objective can be defined to enforce invariance on the embeddings  $z_t$  or representations  $c_t$  of the RPL network. Mathematically, the objective functions corresponding to these variations can be written as:

$$\mathcal{L}_{\text{IL}^z} = \|z_t - \text{SG}(z_{t+\Delta})\|^2 + \mathcal{L}_{\text{reg}}(z_t), \quad (4)$$

and

$$\mathcal{L}_{\text{IL}^c} = \|c_t - \text{SG}(c_{t+\Delta})\|^2 + \mathcal{L}_{\text{reg}}(c_t), \quad (5)$$

respectively with the target for  $L_{\text{reg}}$  changed in each case as indicated.

## Context-predictive learning (CL)

Finally, the RPL objective is not the only possible implementation of a JEPA given our three-component network architecture. Alternatively, one could also optimize the same network for predicting future representations  $c_{t+\Delta}$  from the integrator instead of  $z_{t+\Delta}$ , yielding the context-predictive learning (CL) objective function:

$$\mathcal{L}_{\text{CL}} = \|f(c_t) - \text{SG}(c_{t+\Delta})\|^2. \quad (6)$$

## Datasets and network models

We investigated RPL's performance on multiple synthetic and real-world video and audio datasets with corresponding deep neural network architectures for the encoder, integrator, and predictor.

## Synthetic moving animal video dataset

We generated the synthetic 32-frame moving animal videos by first generating sequences of the continuous latent variables. To that end, we randomly sampled an animal category and simulated its motion with a random initial velocity  $v_1, v_2 \sim U(-8, 8)$  and rotation with a constant randomly sampled angular velocity  $\omega \sim \mathcal{U}(-\frac{\pi}{6}, \frac{\pi}{6})$ . We also simulated motion in a third depth dimension by changing the scale of the animal starting with a random initial velocity  $v_3 \sim U(-0.125, 0.125)$ . Furthermore, we assumed elastic collisions at the position boundaries in all three dimensions, which changed the direction of motion. Overall, the latent variables included the animal category  $S$ , position  $(r_1, r_2)$ , scale  $r_3$ , velocity  $\mathbf{v} = (v_1, v_2, v_3)$ , planar orientation  $\theta$ , and the corresponding angular velocity  $\omega$ . We sampled the initial orientation from a uniform distribution  $\theta \sim U(0, 2\pi)$ . The initial position in the  $r_1$ - $r_2$  plane was sampled from a uniform distribution:  $r_1, r_2 \sim U(8, 56)$ , where 8 is the padding around the image, ensuring that the animal remained visible within the frame. The scale  $r_3$  was randomly initialized from and remained limited to the interval  $[0.2, 1]$ .

We then sampled the latent variable trajectory at discrete intervals and constructed the corresponding  $64 \times 64$  RGB video frames. Each frame contained one antialiased animal image scaled and rendered on a gray background. Here, a scale  $r_3 = 1$  corresponded to an image size of  $32 \times 32$  pixels. Finally, we added independent Gaussian noise  $\sim N(0, 0.1)$  to each pixel and color channel (Supplementary Fig. S1 [↗](#)). Using this procedure, we built a fixed corpus of 4000 32-frame videos for validation and testing, and randomly generated training sets of 16000 videos for each network model we trained.

## Network architecture and training

We used a convolutional neural network (CNN) with six convolutional layers for the encoder, each with 32 filters of size  $5 \times 5$  and padding of size 2, followed by a batch normalization layer and a rectified linear unit (ReLU) activation function. We used a stride of two for the first two layers to reduce the spatial dimensions, and stride one for later layers. The output of the last convolutional layer was flattened to obtain the embedding vector  $\mathbf{z}_t$ . This embedding was then fed to the integrator with a single layer of 512 long short-term memory (LSTM) units to obtain the internal representation  $\mathbf{c}_t$ . Finally, this representation was passed through the predictor, a feedforward Multi-Layer Perceptron (MLP) network with one hidden layer of 512 units and the ReLU activation function, followed by a linear output layer. We trained the networks using the AdamW optimizer [104] with learning rate  $3 \times 10^{-4}$ , weight decay  $10^{-5}$ , and batch size 128 for 1000 epochs using the aforementioned objective functions with  $\Delta = 1$  on the 16000 training videos, i.e., a total of 512000 frames, using a cosine learning rate schedule. In the case of RPL and CL (Eqs. (2) [↗](#) and (6)), we used a learning rate 10 times higher for the predictor network relative to the encoder, an established protocol for JEPA training [37]. For IL variants, we used  $\lambda_1 = 20$  and  $\lambda_2 = 200$  for the regularization strengths.

To compare our model to established SSL, albeit not necessarily biologically plausible algorithms, we trained the same network architecture via an InfoNCE loss [34] with 10 negative samples per positive sample and using either a nonlinear or linear predictor network (Table S1 [↗](#)).

## Evaluation

We then evaluated the trained network models using linear probing as follows. First, we froze all network parameters and computed the internal representations of the training and testing datasets. We then standardized these representations to zero mean and unit variance. We used the resulting representations of the training dataset to fit an ordinary least square linear regression model for each continuous latent variable and a linear classifier for animal category. Since orientation is not a linear variable, we used  $a = \sin \theta$  and  $b = \cos \theta$  as target variables for the linear readout, which was trained by minimizing the mean-squared-error. We additionally optimized a phase offset for each animal category to correct for orientation offset biases per animal in the shared representation space. Moreover, we decoded the angular speed  $|\omega|$  from time-averaged video representations. We trained the linear classifier on time-averaged video frames of the

training set using the AdamW optimizer with a learning rate of  $10^{-3}$ , weight decay of  $10^{-5}$ , and a batch size of 128 for 250 epochs. Throughout the article, we reported the classification accuracies and coefficient of determination  $R^2$  for regression on the held-out test data. Finally, to provide a lower reference value for comparison, we evaluated the readout performance as described above for an untrained, randomly initialized network as a control. Similarly, to compare with supervised learning, we also trained the same network end-to-end using the regression and classification readouts as supervised training objectives for optimizing the encoder and integrator. Note that the predictor is unused during evaluation in all cases, and also during supervised training.

## Internal simulation of trajectories

To investigate the model's ability to accurately simulate trajectories solely based on internal representations, we primed the network with the first eight frames of a video. After the initial priming phase, the predicted embeddings were fed back to the integrator autoregressively for the following sixteen time steps thereby generating simulated internal representations (Fig. 3a). To evaluate the quality of the simulation, we then estimated the latent variables from the internal representations using the linear decoders described above and compared them to the actual latent variables of the ground truth trajectory (cf. Fig. 3 and S7).

## Synthetic sequences of handwritten digits

To test RPL on more abstract stimulus sequences, we generated synthetic sequences of handwritten digit images consisting of digit triplets separated by zeros (Fig. 4a). We generated the digit sequences by sampling clusters, triplets and digits for 64 time steps according to a structured grammar (cf. Fig. 4b). Specifically, we clustered the digits into four groups: {1, 2, 3}, {4, 5, 6}, {7, 8, 9}, and {0}. The permutations of the digits within each of the first three groups constituted the triplets, making a total of 18 possible triplets. Each digit sequence consisted of such triplets interspersed with zeros, starting with a triplet drawn from a randomly chosen cluster. Each triplet was either followed by another random triplet from the same cluster with a probability of 0.8, or a zero with a probability of 0.2. After a zero, a cluster group was chosen uniformly at random to continue the sequence. Overall, the latent variables in this case were the digit identity, the cluster group, and the triplet identity. Given the generated sequence of digits, we finally sampled corresponding  $28 \times 28$  grayscale images from the MNIST training and test sets for training and testing respectively. Using this approach, we generated 10000 training sequences and 10000 test sequences for each model.

## Network architecture and training

For the encoder, we used an MLP with two hidden layers consisting of 32 units, each followed by batch normalization and ReLU activation functions. We fed the resulting embeddings into a single-layer recurrent neural network (RNN) integrator with 128 units to generate the internal representations. The representations were further passed through a feedforward MLP predictor network with a single hidden layer of 128 units and ReLU activation function, followed by a linear output layer. We trained all networks using the AdamW optimizer with a learning rate  $3 \times 10^{-4}$ , weight decay  $10^{-3}$ , and batch size 64 for 200 epochs using the different objective functions with  $\Delta = 1$  on the training sequences. For IL, we used  $\lambda_1 = 1$  and  $\lambda_2 = 10$  for the regularization strength.

## Evaluation

To assess the quality of the learned representations, we first froze the parameters of the trained networks. Then, we extracted the representations of the training and test sets and standardized them as before. We subsequently trained a set of linear classifiers on these representations to decode digit identity and cluster identity at each time step, and triplet identity after observing the first two digits of any triplet subsequence (cf. Fig. 4b). We trained all three linear classifiers for 50 epochs using AdamW with a learning rate of  $10^{-3}$ , weight decay of  $10^{-5}$ , and a batch size of 128. As before, we also computed performance estimates for a random untrained baseline network as well as end-to-end supervised training.

## Real-world videos of behaving mice

To study RPL's ability to extract latent variables from real-world data, we used videos from the open source dataset of behaving mice in an open-field arena provided by Luxem et al. [47]. Concretely, we used two one-hour videos with a 25 Hz frame rate resulting in a total of  $2 \times 90000$  frames. To split the data into training and test sets, we sliced the videos into one-minute segments from which we used the first fifty seconds for training and the last 10 seconds for testing. We downsampled all frames to  $96 \times 96$  pixel frames using HandBrake [105] and subsampled the video segments at 5 Hz, i.e., only using every fifth frame, to decrease computational training cost. We further computed the proxy latent variables velocity, orientation, and position based on the dataset's DeepLabCut [48] estimates of tail root and nose position. To reduce label noise in the velocity estimates of each sample (Supplementary Fig. S9 [49]), we smoothed them using a cumulative summation over time before standardizing the values. We further noted that the dataset contained erroneous nose position outside of the arena, for instance, whenever the nose was not visible when the animal was rearing. To prevent this label noise from entering our orientation estimates, we replaced the affected data points with linearly interpolated values.

## Network architecture and training

We used a CNN with six convolutional layers, each with 32 filters as the encoder. We used  $5 \times 5$  filters with padding and strides of size two for the first three layers, and  $3 \times 3$  filters with padding and strides of size one for the following three layers. All convolutional layers were followed by a batch normalization layer and ReLU activation functions. The output of the last layer was flattened to obtain the embedding vector  $\mathbf{z}_t$ . This embedding was then fed to the integrator consisting of a single layer of 256 LSTM units to obtain the internal representation  $\mathbf{c}_t$ . Finally, the internal representations passed through an MLP predictor network with a single hidden layer comprised of 256 units followed by ReLU activation functions and a dense linear output layer. We trained the networks using the AdamW optimizer with learning rate  $3 \times 10^{-4}$ , weight decay  $10^{-5}$ , and batch size 128 for 500 epochs using the aforementioned objective functions with  $\Delta = 1$  on the training set. In the case of RPL (Eq. (2) [49]), we trained the predictor network with a 10 times faster learning rate than the rest of the network. When training using an IL objective, we additionally set the regularization strengths to  $\lambda_1 = 20$  and  $\lambda_2 = 200$ .

## Evaluation

We evaluated the networks on the test subset of the dataset using linear least squares with L2 regularization, i.e., ridge regression, with regularization coefficient of 1.0 on the trained internal representations to decode the latent variables  $\mathbf{r}_N$ ,  $\mathbf{v}_T$ , and  $\theta$ .

## Real-world natural speech recordings

To test our representation learning framework in a different data domain, we used a 100-hour audio recording subset of the Librispeech dataset [49] for training and evaluating the network. Librispeech contains English audio book files from 251 speakers, with 41 force-aligned phoneme labels. The audio files were sampled at 16 kHz and there is a phoneme label for roughly every 10 ms of audio.

## Network architecture and training

Following the work of Oord et al. [34], we used a 1D CNN architecture acting on the raw audio waveform with 512 output units as the encoder. The resulting embeddings were then fed into the integrator, a single LSTM layer with 256 units, to generate the internal representations. The representations were then passed through an MLP predictor network with a single hidden layer of 512 units and ReLU activation function, followed by a dense linear output layer. The network was trained using the AdamW optimizer with learning rate  $3 \times 10^{-4}$ , no weight decay, and batch size 256

for 500 epochs using the different objective functions with  $\Delta = 8$ . The hyperparameter  $\Delta$  was determined via grid search to be the optimal prediction horizon for RPL. In the case of IL objectives, we used  $\lambda_1 = 2$  and  $\lambda_2 = 20$  as regularization strengths.

## Evaluation

We evaluated the trained networks on the test set following the same steps as above and by training linear readouts to decode the phoneme and speaker labels. The linear readout contained 251 units for speaker classification and 41 units for phoneme classification. We trained the linear classifiers using the AdamW optimizer with learning rate  $10^{-3}$ , weight decay  $10^{-5}$ , and batch size 64 (for speaker identity) for 50 epochs. Following the same procedure as in Oord et al. [34], we use a batch size of one for training readouts for the phoneme classification task, as the start and end point annotations for phoneme labels are not provided. As before, we evaluated a randomly initialized network for comparison.

## Learning successor-like representations

To investigate whether RPL learns successor-like representations, we modeled a stimulus sequence A-B-C-D as in the experiment [45] with one-hot vectors. Since the one-hot encoding of the inputs alleviated the need for learning stimulus representations, we replaced the encoder with an identity matrix. For simplicity, we used a vanilla RNN with four units as integrator, whose weights we initialized with an identity matrix, and a dense linear predictor whose weights we initialized with a Glorot uniform distribution [106]. We trained  $n = 35$  networks for 352 epochs on the A-B-C-D sequence using the RPL objective with a learning rate of  $10^{-3}$  for the integrator and 0.01 for the predictor. Further, we used a weight decay of  $10^{-4}$  for all the components.

We trained logistic regression models on the internal representations of the full sequences to classify the sequence locations. We then sorted the integrator units based on their selectivity to sequence locations. We computed the average weight matrix of all logistic regression models and visualized it to see if it resembles a successor matrix (Fig. 6b). To generate the partial sequences, we used zero input vectors for the non-stimulated sequence locations and the corresponding one-hot vector for the stimulated location. We then computed the activity at stimulated and non-stimulated locations for each partial sequence by recording the mean population activity of the integrator units at each sequence location and subtracting the population response to a zero input (Fig. 6d). These values were averaged over all partial sequences for each stimulated and non-stimulated location to obtain the activity at individual predecessor, stimulated, and successor locations (Fig. 6c).

## Learning abstract sequence representations from local-global oddballs

To study whether RPL learns abstract sequence representations similar to the primate PFC, we modeled the local-global oddball experiment conducted by Bellet et al. [44]. We first extracted images of the ten visual stimuli used in the experiment and downsampled them to  $32 \times 32$  pixel color images. We then constructed the oddball sequences from [44] consisting of pairs of stimuli organized in blocks of twenty sequences (scaled down from the original block size of two hundred). Each sequence comprised four stimulus presentations, either repeating the same stimulus (xxxx or xx for short) or with a local deviant at the last position (xxxY or xY for short). Here, x stands for either of the two items in each pair of stimuli and y for the other, fixed within one block. The first five sequences in each block were consistently either xx or xY, setting the global standard or “context” for that block. The next fifteen sequences were randomly set to either global standards (80%) or rare global deviants (20%). The global standards were repeats of the standard sequence, i.e., xx in an xx block (xx | xx) or xY in an xY block (xY | xY). The global deviants differed from the standard sequence only in the last position (xY | xx or xx | xY). The full sequence of stimuli also contained blank frames, one between each pair of consecutive stimuli in a given

sequence, and four between sequences in the same block. We generated 2000 such blocks with the stimulus pairs and contexts uniformly chosen at random for each block for training the network model, and 2000 more for testing.

## Network architecture and training

Since the stimulus set consisted of only ten images, we used a simple MLP with two hidden layers and batch normalization for the encoder. The MLP took in the full vector of image pixels as input and had 64 units in each hidden layer and its output layer. We then fed the embedding from the output of the encoder to an integrator consisting of a single LSTM layer with 256 units to generate the internal representations. These were then passed through an MLP predictor with one hidden layer with 256 units. Both the encoder and predictor MLPs used the ReLU activation function. For training, each block of 20 sequence presentations was given as one training sample sequence. We trained the network using the AdamW optimizer with learning rate  $3 \times 10^{-4}$ , batch size 128, and weight decay  $10^{-3}$  for 2000 epochs using the RPL objective function with  $\Delta = 1$ . As before, we also used a 10 times higher learning rate for the predictor.

## Evaluation

For evaluating the trained network, we closely followed the analysis of PFC population recordings done by Bellet et al. [44]. This consisted of training linear decoders for serial sequence position, stimulus identity, global context and presence of local and global deviance as above, but with specific sequence subtypes used for training the readout (cf. Fig 3 and Methods in [44]) as summarized below:

### Serial position

For sequence position, we trained binary classifiers using logistic regression for each position in the sequence using only xx trials for training. This yielded predicted probabilities of items 1 through 4 at each time point of a sequence, plotted for the held-out test set in Fig. 7b.

### Stimulus identity

Stimulus identity was defined as a fixed binary label for the two items in each of the five stimulus pairs. We trained a binary classifier for decoding this binary label using only the first three stimuli in each sequence for training. The projections of the representations onto the weights found by the binary classifier are plotted in Fig. 7c for the held-out test set of sequences.

### Global context

For decoding the context, we first averaged the representations over the first three stimulus pairs in every sequence attaching a binary label corresponding to whether the global standard was a repeat sequence (xx) or a sequence with a local oddball (xY). As above, we trained a binary classifier with logistic regression and plotted the projections of the representations from held-out sequences onto the weights found by training the binary classifier.

### Deviance

For decoding both global and local deviance, we used the representation corresponding to the last stimulus presentation to train a binary classifier with logistic regression. Then, we plotted the projections of the representations for the test set onto the weights of the binary classifier.

All linear decoders were trained with AdamW using the learning rate  $10^{-3}$  and weight decay  $10^{-5}$ .

### Hierarchical RPL (h-RPL)

To assess the effectiveness of RPL in a hierarchical setting, we constructed a network from six RPL circuits each with a shallow single-layer 2D CNN with 32 filters of size  $5 \times 5$  and padding of size 2, followed by a batch normalization layer and a ReLU activation function. The first two convolutional layers had a stride of two while the convolutional layers in higher circuits had a

stride of 1. Note that this stack of shallow instantaneous encoders in h-RPL results in the same deep convolutional encoder used in RPL for the same dataset. The encoder output was sent to the next circuit in the hierarchy and a flattened version was sent to the local integrator circuit, an LSTM with 512 units (Fig. 8a) for all circuits. Each circuit also had its own MLP predictor with a single hidden layer comprised of 512 units with ReLU activation functions and followed by a dense linear output. Importantly, the different circuits were gradient isolated and independently optimized the RPL objective (Eq. 2) at each level of the hierarchy, without backpropagating the gradients between circuits. Finally, we trained the model on the synthetic moving animal videos (cf. Fig. 2) using the AdamW optimizer with learning rate  $3 \times 10^{-4}$ , weight decay  $10^{-5}$ , and batch size 128 for 500 epochs using the RPL objective. As before, we evaluated the representations using linear probing albeit now for each circuit individually.

## Data and code availability

The code to reproduce all modeling results in this article is available on Github [110]. We further used the following publicly available datasets: The Librispeech dataset is freely available for download [107]. We used train and test splits and force-aligned phoneme annotations as defined previously [34] and made available by the authors (<https://drive.google.com/drive/folders/1BhJ2umKH3whguxMwifaKtSra0TgAbtfb>). The mouse arena videos are available on Figshare [108]. Further, the two down-sampled videos used in this article are included in our Github repository. Finally, the plots for the macaque PFC experiment [44] were reproduced using data [109] and analysis code [111] shared publicly by the authors.

## Acknowledgements

We thank Anna Vasilevskaya, Georg Keller, and all Zenke Lab members for inspiring discussions throughout this project. Special thanks to Peter Buttaroni for contributing illustrations for figures. This project was supported by the Swiss National Science Foundation (Grant Number PCEFP3 202981), EU's Horizon Europe Research and Innovation Programme (Grant Agreement No. 101070374, CONVOLVE) funded through SERI (Ref. 1131–52302), and the Novartis Research Foundation.

## Additional information

### Author contributions

Ashena Gorgan Mohammadi & Manu Srinath Halvagal: Formal analysis, Investigation, Methodology, Software, Visualization, Writing – original draft, Writing – review and editing; Friedemann Zenke: Conceptualization, Funding acquisition, Methodology, Supervision, Writing – original draft, Writing – review and editing.

### Funding

Funder	Grant reference number	Author
Swiss National Science Foundation	PCEFP3 202981	Friedemann Zenke
Horizon Europe Research Innovation Program	101070374	Friedemann Zenke
Novartis Research Foundation		Friedemann Zenke

### Author ORCID iDs

Ashena Gorgan Mohammadi: <https://orcid.org/0000-0002-8998-3191>

Manu Srinath Halvagal: <https://orcid.org/0000-0001-8374-8257>

Friedemann Zenke: <https://orcid.org/0000-0003-1883-644X>

## Additional files

[Supplementary Material](#)

## References

- [1] Logothetis Nikos K, et al. (1996) Visual object recognition. *Annual review of neuroscience* **19**:577-621 <https://doi.org/10.1146/annurev.ne.19.030196.003045> | PubMed
- [2] Hubel David H, et al. (1962) Receptive fields, binocular interaction and functional architecture in the cat's visual cortex. *The Journal of physiology* **160**:106 <https://doi.org/10.1113/jphysiol.1962.sp006837> | PubMed
- [3] Duhamel Jean-René, et al. (1992) The updating of the representation of visual space in parietal cortex by intended eye movements. *Science* 90-92 <https://doi.org/10.1126/science.1553535> | PubMed
- [4] Whitney David, et al. (2003) Flexible retinotopy: motion-dependent position coding in the visual cortex. *Science* **302**:878-881 <https://doi.org/10.1126/science.1087839> | PubMed
- [5] Johnson Philippa Anne, et al., et al. (2023) Position representations of moving objects align with real-time position in the early visual response. *eLife* **12**:e82424 <https://doi.org/10.7554/eLife.82424> | PubMed
- [6] Maunsell John H, et al. (1983) Functional properties of neurons in middle temporal visual area of the macaque monkey. I. Selectivity for stimulus direction, speed, and orientation. *Journal of neurophysiology* **49**:1127-1147 <https://doi.org/10.1152/jn.1983.49.5.1127> | PubMed
- [7] Hennestad Eivind, et al. (2021) Mapping vestibular and visual contributions to angular head velocity tuning in the cortex. *Cell Reports* **37** <https://doi.org/10.1016/j.celrep.2021.110134> | PubMed
- [8] Priebe Nicholas J., et al. (2006) Tuning for Spatiotemporal Frequency and Speed in Directionally Selective Neurons of Macaque Striate Cortex. *Journal of Neuroscience* **26**:2941-2950 <https://doi.org/10.1523/JNEUROSCI.3936-05.2006> | PubMed
- [9] Roitman Alexander V., et al. (2005) Position, Direction of Movement, and Speed Tuning of Cerebellar Purkinje Cells during Circular Manual Tracking in Monkey. *Journal of Neuroscience* **25**:9244-9257 <https://doi.org/10.1523/jneurosci.1886-05.2005> | PubMed
- [10] Illing Bernd, et al. (2021) Local plasticity rules can learn deep representations using self-supervised contrastive predictions. *Advances in Neural Information Processing Systems* **34**
- [11] Halvagal Manu Srinath, et al. (2023) The combination of Hebbian and predictive plasticity learns invariant object representations in deep sensory networks. *Nat Neurosci* **26**:1906-1915 <https://doi.org/10.1038/s41593-023-01460-y> | PubMed
- [12] Keller Georg B., et al. (2018) Predictive Processing: A Canonical Cortical Computation. *Neuron* **100**:424-435 <https://doi.org/10.1016/j.neuron.2018.10.003> | PubMed
- [13] Rao R P, et al. (1999) Predictive coding in the visual cortex: A functional interpretation of some extra-classical receptive-field effects. *Nature Neuroscience* **2**:79-87 <https://doi.org/10.1038/4580> | PubMed
- [14] Singer Yosef, et al. (2018) Sensory cortex is optimized for prediction of future input. *eLife* **7**:e31557 <https://doi.org/10.7554/eLife.31557> | PubMed
- [15] Lotter William, et al. (2020) A neural network trained for prediction mimics diverse features of biological neurons and perception. *Nature machine intelligence* **2** <https://doi.org/10.1038/s42256-020-0170-9> | PubMed
- [16] Hertäg Loreen, et al. (2020) Learning prediction error neurons in a canonical interneuron circuit. *eLife* **9**:e57541 <https://doi.org/10.7554/eLife.57541> | PubMed
- [17] Recanatesi Stefano, et al. (2021) Predictive learning as a network mechanism for extracting low-dimensional latent space representations. *Nature Communications* **12**:1417 <https://doi.org/10.1038/s41467-021-21696-1> | PubMed
- [18] Mikulasch Fabian A., et al. (2023) Where is the error? Hierarchical predictive coding through dendritic error computation. *Trends in Neurosciences* **46**:45-59 <https://doi.org/10.1016/j.tins.2022.09.007> | PubMed

- [19] Levenstein Daniel, et al. (2024) Sequential predictive learning is a unifying theory for hippocampal representation and replay. *bioRxiv* <https://doi.org/10.1101/2024.04.28.591528>
- [20] Nejad Kevin Kermani, et al. (2025) Self-supervised predictive learning accounts for cortical layer-specificity. *Nat Commun* **16**:6178 <https://doi.org/10.1038/s41467-025-61399-5> | PubMed
- [21] Rao Rajesh PN (2024) A sensory–motor theory of the neocortex. *Nature neuroscience* **27**:1221-1235 <https://doi.org/10.1038/s41593-024-01673-9> | PubMed
- [22] Taylor Luke, et al. (2024) Temporal prediction captures key differences between spiking excitatory and inhibitory V1 neurons. *bioRxiv* 2024.05.12.593763 <https://doi.org/10.1101/2024.05.12.593763>
- [23] Tscshantz Alexander, et al. (2023) Hybrid predictive coding: Inferring, fast and slow. *PLOS Computational Biology* **19**:e1011280 <https://doi.org/10.1371/journal.pcbi.1011280> | PubMed
- [24] Oliviers Gaspard, et al. (2025) Bidirectional predictive coding. *arXiv* <https://doi.org/10.48550/arXiv.2505.23415>
- [25] Westerberg Jacob A., et al. (2025) Hierarchical interactions between sensory cortices defy predictive coding. *Trends in Cognitive Sciences* <https://doi.org/10.1016/j.tics.2025.09.018> | PubMed
- [26] Solomon Selina S, et al. (2021) Limited Evidence for Sensory Prediction Error Responses in Visual Cortex of Macaques and Humans. *Cerebral Cortex* **31**:3136-3152 <https://doi.org/10.1093/cercor/bhab014> | PubMed
- [27] Gillon Colleen J., et al. (2024) Responses to Pattern-Violating Visual Stimuli Evolve Differently Over Days in Somata and Distal Apical Dendrites. *J. Neurosci* **44** <https://doi.org/10.1523/JNEUROSCI.1009-23.2023> | PubMed
- [28] Richter David, et al. (2024) High-Level Visual Prediction Errors in Early Visual Cortex. *PLOS Biology* **22**:e3002829 <https://doi.org/10.1371/journal.pbio.3002829> | PubMed
- [29] Furutachi Shohei, et al. (2024) Cooperative thalamocortical circuit mechanism for sensory prediction errors. *Nature* **633**:398-406 <https://doi.org/10.1038/s41586-024-07851-w> | PubMed
- [30] Vasilevskaya Anna, et al. (2026) A functional influence based circuit motif that constrains the set of plausible algorithms of cortical function. *bioRxiv* <https://doi.org/10.64898/2026.01.29.702557>
- [31] Balestriero Randall, et al. (2023) A Cookbook of Self-Supervised Learning. *arXiv* 2304.12210 <https://doi.org/10.48550/arxiv.2304.12210>
- [32] Bardes Adrien, et al. (2022) VICReg: Variance-Invariance-Covariance Regularization for Self-Supervised Learning. In: International Conference on Learning Representations.
- [33] Niu Julie Xueyan, et al. (2024) Learning predictable and robust neural representations by straightening image sequences. *Advances in Neural Information Processing Systems* **37**:40316-40335 <https://doi.org/10.52202/079017-1275>
- [34] van den Oord Aaron, et al. (2018) Representation Learning with Contrastive Predictive Coding. *arXiv* 1807.03748 <https://doi.org/10.48550/arXiv.1807.03748>
- [35] LeCun Yann (2022) A Path Towards Autonomous Machine Intelligence Version 0.9.2, 2022-06-27. *OpenReview.net*
- [36] Grill Jean-Bastien, et al. (2020) Bootstrap Your Own Latent: A New Approach to Self-Supervised Learning. *arXiv* 2006.07733 <https://doi.org/10.48550/arxiv.2006.07733>
- [37] Chen Xinlei, et al. (2021) Exploring simple siamese representation learning. In: *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 15750-15758 <https://doi.org/10.1109/cvpr46437.2021.01549>
- [38] Assran Mahmoud, et al. (2023) Self-supervised learning from images with a joint-embedding predictive architecture. In: *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 15619-15629 <https://doi.org/10.1109/cvpr52729.2023.01499>

- [39] Ghaemi Hafez, et al. (2024) Seq-JEPA: Autoregressive Predictive Learning of Invariant-Equivariant World Models. In: *NeurIPS 2024 Workshop: Self-Supervised Learning - Theory and Practice*.
- [40] Bardes Adrien, et al. (2024) V-JEPA: Latent Video Prediction for Visual Representation Learning. *OpenReview.net*
- [41] Assran Mido, et al. (2025) V-jepa 2: Self-supervised video models enable understanding, prediction and planning. *arXiv* 2506.09985 <https://doi.org/10.48550/arxiv.2506.09985>
- [42] Balestriero Randall, et al. (2024) Learning by Reconstruction Produces Uninformative Features For Perception. *arXiv* 2402.11337 <https://doi.org/10.48550/arXiv.2402.11337>
- [43] Littwin Etai, et al. (2024) How jepa avoids noisy features: The implicit bias of deep linear self distillation networks. *arXiv* 2407.03475 <https://doi.org/10.48550/arxiv.2407.03475>
- [44] Bellet Marie E, et al. (2024) Spontaneously emerging internal models of visual sequences combine abstract and event-specific information in the prefrontal cortex. *Cell reports* **43** <https://doi.org/10.1016/j.celrep.2024.113952> | [PubMed](#)
- [45] Ekman Matthias, et al. (2023) Successor-like representation guides the prediction of future events in human visual cortex and hippocampus. *eLife* **12**:e78904 <https://doi.org/10.7554/eLife.78904> | [PubMed](#)
- [46] Deng Li (2012) The mnist database of handwritten digit images for machine learning research. *IEEE Signal Processing Magazine* **29**:141-142 <https://doi.org/10.1109/msp.2012.2211477>
- [47] Luxem Kevin, et al. (2022) Identifying behavioral structure from deep variational embeddings of animal motion. *Communications Biology* **5** <https://doi.org/10.1038/s42003-022-04080-7> | [PubMed](#)
- [48] Mathis Alexander, et al. (2018) DeepLabCut: markerless pose estimation of user-defined body parts with deep learning. *Nature Neuroscience* **21** <https://doi.org/10.1038/s41593-018-0209-y> | [PubMed](#)
- [49] Panayotov Vassil, et al. (2015) Librispeech: An ASR corpus based on public domain audio books. In: *2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)* 2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP). pp. 5206-5210 <https://doi.org/10.1109/ICASSP.2015.7178964>
- [50] Dayan Peter (1993) Improving generalization for temporal difference learning: The successor representation. *Neural computation* **5**:613-624 <https://doi.org/10.1162/neco.1993.5.4.613>
- [51] Lillicrap Timothy P., et al. (2020) Backpropagation and the brain. *Nature Reviews Neuroscience* **23**:335-346 <https://doi.org/10.1038/s41583-020-0277-3> | [PubMed](#)
- [52] Cajal Santiago Ramon y (1899) *Comparative study of the sensory areas of the human cortex* Wikimedia Commons. Public Domain.
- [53] Lillicrap Timothy P, et al. (2016) Random synaptic feedback weights support error backpropagation for deep learning. *Nature communications* **7** <https://doi.org/10.1038/ncomms13276> | [PubMed](#)
- [54] Nøkland Arild (2016) Direct feedback alignment provides learning in deep neural networks. *Advances in neural information processing systems* **29**
- [55] Whittington James CR, et al. (2017) An approximation of the error backpropagation algorithm in a predictive coding network with local hebbian synaptic plasticity. *Neural computation* **29**:1229-1262 [https://doi.org/10.1162/neco\\_a\\_00949](https://doi.org/10.1162/neco_a_00949) | [PubMed](#)
- [56] Sacramento João, et al. (2018) Dendritic cortical microcircuits approximate the backpropagation algorithm. *Advances in neural information processing systems* **31**
- [57] Hénaff Olivier J, et al. (2019) Perceptual straightening of natural videos. *Nature neuroscience* **22**:984-991 <https://doi.org/10.1038/s41593-019-0377-4> | [PubMed](#)
- [58] Hénaff Olivier J, et al. (2021) Primary visual cortex straightens natural video trajectories. *Nature communications* **12** <https://doi.org/10.1038/s41467-021-25939-z> | [PubMed](#)
- [59] Bastos Andre M., et al. (2012) Canonical Microcircuits for Predictive Coding. *Neuron* **76**:695-711 <https://doi.org/10.1016/j.neuron.2012.10.038> | [PubMed](#)
- [60] Hawkins Jeff (2021) *A Thousand Brains: A New Theory of Intelligence* Hachette UK.

- [61] Brunec Iva K., et al. (2022) Predictive Representations in Hippocampal and Prefrontal Hierarchies. *Journal of Neuroscience* **42**:299-312 <https://doi.org/10.1523/JNEUROSCI.1327-21.2021> | PubMed
- [62] Price Byron H, et al. (2023) Expectation violations produce error signals in mouse V1. *Cerebral Cortex* **33**:8803-8820 <https://doi.org/10.1093/cercor/bhad163> | PubMed
- [63] Tang Matthew F, et al. (2023) Expectation violations enhance neuronal encoding of sensory information in mouse primary visual cortex. *Nature Communications* **14** <https://doi.org/10.1038/s41467-023-36608-8> | PubMed
- [64] Han Shuting, et al. (2024) Behavior-Relevant Top-down Cross-Modal Predictions in Mouse Neocortex. *Nature Neuroscience* 1-11 <https://doi.org/10.1038/s41593-023-01534-x> | PubMed
- [65] Boeshertz Gauthier, et al. (2025) Predictive learning enables compositional representations. *bioRxiv* <https://doi.org/10.1101/2025.09.26.678731>
- [66] Van Assel Hugues, et al. (2025) Joint Embedding vs Reconstruction: Provable Benefits of Latent Space Prediction for Self Supervised Learning. *arXiv* 2505.12477 <https://doi.org/10.48550/arxiv.2505.12477>
- [67] Zmarz Pawel, et al. (2016) Mismatch Receptive Fields in Mouse Visual Cortex. *Neuron* **92**:766-772 <https://doi.org/10.1016/j.neuron.2016.09.057> | PubMed
- [68] Ed Li T., et al. (2025) Predictive Coding Model Detects Novelty on Different Levels of Representation Hierarchy. *Neural Computation* **37**:1373-1408 [https://doi.org/10.1162/neco\\_a\\_01769](https://doi.org/10.1162/neco_a_01769) | PubMed
- [69] Ororbia Alexander, et al. (2025) Meta-Representational Predictive Coding: Biomimetic Self-Supervised Learning. *arXiv* 2503.21796 [https://doi.org/10.31219/osf.io/s7mcd\\_v1](https://doi.org/10.31219/osf.io/s7mcd_v1)
- [70] Delrocq Ariane, et al. (2024) Critical periods support representation learning in a model of cortical processing. *bioRxiv* <https://doi.org/10.1101/2024.12.20.629674>
- [71] Wiskott Laurenz, et al. (2002) Slow Feature Analysis: Unsupervised Learning of Invariances. *Neural Computation* **14**:715-770 <https://doi.org/10.1162/089976602317318938> | PubMed
- [72] Asabuki Toshitake, et al. (2025) Learning predictive signals within a local recurrent circuit. *Proceedings of the National Academy of Sciences* **122**:e2414674122 <https://doi.org/10.1073/pnas.2414674122> | PubMed
- [73] Noroozi Mehdi, et al. (2016) Unsupervised Learning of Visual Representations by Solving Jigsaw Puzzles. In: *Lecture Notes in Computer Science Computer Vision – ECCV 2016*. pp. 69-84 [https://doi.org/10.1007/978-3-319-46466-4\\_5](https://doi.org/10.1007/978-3-319-46466-4_5)
- [74] Zbontar Jure, et al. (2021) Barlow Twins: Self-Supervised Learning via Redundancy Reduction. In: *arXiv Proceedings of the 38th International Conference on Machine Learning*. **139** Proceedings of Machine Learning Research. PMLR. pp. 12310-12320
- [75] Richthofer Stefan, et al. (2015) Predictable feature analysis. In: 2015 IEEE 14th international Conference on machine learning and applications (ICMLA). IEEE. pp. 190-196 <https://doi.org/10.1109/icmla.2015.158>
- [76] Vladimirskiy Boris, et al. (2015) Hierarchical Novelty-Familiarity Representation in the Visual System by Modular Predictive Coding. *PLOS ONE* **10**:1-19 <https://doi.org/10.1371/journal.pone.0144636> | PubMed
- [77] Devillers Alexandre, et al. (2022) Equimod: An equivariance module to improve self-supervised learning. *arXiv* 2211.01244 <https://doi.org/10.48550/arxiv.2211.01244>
- [78] Park Jung Yeon, et al. (2022) Learning symmetric embeddings for equivariant world models. *arXiv* 2204.11371 <https://doi.org/10.48550/arxiv.2204.11371>
- [79] Garrido Quentin (2023) Self-supervised learning of split invariant equivariant representations. Université Gustave Eiffel. <https://doi.org/10.70675/ac2d4d52z596fz4742z8a22zbc1504a44a90>
- [80] Wang Yifei, et al. (2024) Understanding the role of equivariance in self-supervised learning. *Advances in Neural Information Processing Systems* **37**:127483-127510 <https://doi.org/10.52202/079017-4048>

- [81] Velarde Osvaldo M., et al. (2024) Recurrent Joint Embedding Predictive Architecture with Recurrent Forward Propagation Learning. *arXiv* <https://doi.org/10.48550/arXiv.2411.16695>
- [82] Histed Mark H. (2025) Active filtering: a predictive function of recurrent circuits of sensory cortex. *arXiv* 2501.10521 <https://doi.org/10.48550/arXiv.2501.10521>
- [83] Xu Shengjin, et al. (2012) Activity recall in a visual cortical ensemble. *Nature neuroscience* **15**:449-455 <https://doi.org/10.1038/nn.3036> | PubMed
- [84] Palmer Stephanie E, et al. (2015) Predictive information in a sensory population. *Proceedings of the National Academy of Sciences* **112**:6908-6913 <https://doi.org/10.1073/pnas.1506855112> | PubMed
- [85] Piasini Eugenio, et al. (2021) Temporal stability of stimulus representation increases along rodent visual cortical hierarchies. *Nature communications* **12** <https://doi.org/10.1038/s41467-021-24456-3> | PubMed
- [86] Murray John D., et al. (2014) A hierarchy of intrinsic timescales across primate cortex. *Nat Neurosci* **17**:1661-1663 <https://doi.org/10.1038/nn.3862> | PubMed
- [87] Stachenfeld Kimberly L., et al. (2017) The hippocampus as a predictive map. *Nat Neurosci* **20**:1643-1653 <https://doi.org/10.1038/nn.4650> | PubMed
- [88] Gershman Samuel J. (2018) The Successor Representation: Its Computational Logic and Neural Substrates. *J. Neurosci* **38**:7193-7200 <https://doi.org/10.1523/JNEUROSCI.0151-18.2018> | PubMed
- [89] Whittington James C. R., et al. (2022) How to build a cognitive map. *Nat Neurosci* **25**:1257-1272 <https://doi.org/10.1038/s41593-022-01153-y> | PubMed
- [90] Lisman J. (1999) Relating hippocampal circuitry to function: recall of memory sequences by reciprocal dentate-CA3 interactions. *Neuron* **22**:233-242 [https://doi.org/10.1016/s0896-6273\(00\)81085-5](https://doi.org/10.1016/s0896-6273(00)81085-5) | PubMed
- [91] Buckner Randy L. (2010) The Role of the Hippocampus in Prediction and Imagination. *Annual Review of Psychology* **61**:27-48 <https://doi.org/10.1146/annurev.psych.60.110707.163508> | PubMed
- [92] Schapiro Anna C., et al. (2017) Complementary learning systems within the hippocampus: a neural network modelling approach to reconciling episodic memory with statistical learning. *Philosophical Transactions of the Royal Society B: Biological Sciences* **372**:20160049 <https://doi.org/10.1098/rstb.2016.0049> | PubMed
- [93] Buzsáki György, et al. (2018) Space and time: the hippocampus as a sequence generator. *Trends in cognitive sciences* **22**:853-869 <https://doi.org/10.1016/j.tics.2018.07.006> | PubMed
- [94] Vezoli Julien, et al. (2021) Cortical hierarchy, dual counterstream architecture and the importance of top-down generative networks. *NeuroImage* **225**:117479 <https://doi.org/10.1016/j.neuroimage.2020.117479> | PubMed
- [95] Markov Nikola T., et al. (2013) Cortical High-Density Counterstream Architectures. *Science* **342**:1238406 PubMed | <https://doi.org/10.1126/science.1238406>
- [96] Guerguiev Jordan, et al. (2017) Towards deep learning with segregated dendrites. *eLife* **6**:e22901 <https://doi.org/10.7554/eLife.22901> | PubMed
- [97] Payeur Alexandre, et al. (2021) Burst-dependent synaptic plasticity can coordinate learning in hierarchical circuits. *Nature Neuroscience* **24**:1010-1019 <https://doi.org/10.1038/s41593-021-00857-x> | PubMed
- [98] Laborieux Axel, et al. (2022) Holomorphic Equilibrium Propagation Computes Exact Gradients Through Finite Size Oscillations. *NeurIPS* **35**:12950-12963 <https://doi.org/10.52202/068431-0941>
- [99] Rossbroich Julian, et al. (2023) Dis-inhibitory neuronal circuits can control the sign of synaptic plasticity. *Advances in Neural Information Processing Systems* **36** <https://doi.org/10.52202/075280-2799>
- [100] Song Yuhang, et al. (2024) Inferring neural activity before plasticity as a foundation for learning beyond backpropagation. *Nature Neuroscience* 348-358 <https://doi.org/10.1038/s41593-023-01514-1> | PubMed

- [101] Goemaere Cédric, et al. (2025) ePC: Overcoming Exponential Signal Decay in Deep Predictive Coding Networks. *arXiv* 2505.20137 <https://doi.org/10.48550/arXiv.2505.20137>
- [102] Tian Yuandong, et al. (2021) Understanding self-supervised learning dynamics without contrastive pairs. In: *arXiv International Conference on Machine Learning*. PMLR. pp. 10268-10278 <https://doi.org/10.48550/arxiv.2102.06810>
- [103] Halvagal Manu Srinath, et al. (2023) Implicit variance regularization in non-contrastive SSL. *Advances in Neural Information Processing Systems* **36**:63409-63436 <https://doi.org/10.52202/075280-2768>
- [104] Loshchilov Ilya, et al. (2019) Decoupled Weight Decay Regularization. *arXiv* 1711.05101 <https://doi.org/10.48550/arXiv.1711.05101>
- [105] Petit Eric (2003) HandBrake: The open source video transcoder. <https://handbrake.fr/>
- [106] Glorot Xavier, et al. (2010) Understanding the difficulty of training deep feedforward neural networks. In: *Proceedings of the thirteenth international conference on artificial intelligence and statistics*. pp. 249-256
- [107] Panayotov V, Chen G, Povey D, Khudanpur S (2015) Librispeech: An ASR corpus based on public domain audio books. *Open Speech and Language Resources*. ID SLR12 <https://www.openslr.org/12>
- [108] Luxem K (2022) VAME Data. figshare. <https://doi.org/10.6084/m9.figshare.19213272>
- [109] Bellet M (2024) Electrophysiological data from the local-global experiment. figshare. <https://doi.org/10.6084/m9.figshare.24311767.v1>
- [110] Mohammadi AG (2025) Recurrent Predictive Learning (RPL). GitHub. <https://github.com/fmi-basel/recurrent-predictive-learning>
- [111] Bellet ME (2024) PFCSequenceModels. GitHub. <https://github.com/mebellet/PFCSequenceModels>

## Peer reviews

### Reviewer #1 (Public review):

Summary:

The paper describes a biologically plausible version of JEPa using recurrent neural networks called RPL for recurrent predictive learning. Given an embedding  $z_t$ , a recurrent neural network processes these inputs with the form:  $c_{t+1} = \text{RNN}(c_t, z_t)$ . Then the predictive network  $f$  is predicting the future inputs with the format:  $\min || f(c_t) - \text{stop\_grad}(z_{t+\Delta t}) ||^2$ . I understand that a prediction error is defined as:  $e = z_{t+\Delta t} - f(c_t)$  to model cortical measurements in the oddball task.

The RPL model is also shown to build an internal world model, with "real-world" data like the movement of moving animals or speech signals. The representation is then compared to V1 data and expected prediction error signals in an oddball setting. In a stacked hierarchy of RNN learning with RPL, the higher layers appear to learn high-level latent variables, although gradients are not propagated downward to the lower layers.

Strengths:

(1) The paper tackles an open question: Self-supervised learning is thought to be a fundamental principle to explain how computation is structured in the brain. Cortical data suggest qualitatively that prediction error is a core principle of representation learning in the brain, but the field is still looking for a simple yet expressive model that would explain how the cortex learns its representations. RPL contributes in that direction by making a useful link between cortical representation learning in RNN models and the JEPa learning algorithm that was demonstrated to scale to large world model learning from video data by Lecun's group. It is very useful to connect this popular deep learning algorithm to cortical data.

(2) The model formalism is relatively elegant and simple: Simple next input prediction objectives are conceptually simple but not necessarily trivial to build at scale. There is a clear benefit in comparison with contrastive or IL methods because they are free from dataset-specific data augmentation and negative samples. Thereby moving the comp neuro field towards conceptually simpler models of representation in the cortex. Yet predictive only models (and in particular predictive models in latent space instead of pixel space) are not easy to build in a stable fashion. JEPA family is basically intended to solve this question; it is very nice and timely to bring this to comp neuro.

(3) The methodology combining comp neuro and deep learning makes sense: The conceptual and qualitative analogy with cortical prediction errors is relevant and consistent with what is expected as a model of self-supervised learning in cortical models. The methodology to compare RPL with IL and CL is methodologically meaningful and grounded: showing, for instance, how some of the models fail to represent some latent structure in some toy datasets is interesting.

(4) h-RPL: The h-RPL is perhaps the most creative departure from the JEPA model family. It would be interesting to say more about what was particularly difficult to see in the latent variables emerging in the hierarchical model. I often find it magical that layer-wise learning rules of this type are not learning redundant representations. Any insights why this is not the case here would be potentially insightful.

Weaknesses:

In general, I fully support the type of question and ideas that the paper is putting forward. It is, however, very hard in this research field to gain insight into specific conceptual contributions or specific bits of experimental data that the model puts forward. In pointing to the following weaknesses, I am encouraging the authors to lay out more clearly what the unique hypothesis is or the contribution of the RPL model that we should remember it for.

(1) The devil is in the details:

1a) Comparison with JEPA variants: JEPA variants are integrating different details into the learning algorithm. Integrating, for instance, "masking" of the latent encoder targets, or EMA in the style of BYOL or Siamese networks, for the predicted representations. It is great that RPL does not seem to need any of those (next input prediction is a natural implementation of masking, and EMA does not seem to be used). It is notoriously hard for the JEPA model to work without these features. Since some of these details are sometimes surprisingly crucial for a simulation to work, it would be good to report which of the other important details were key to live without EMA and masking. Is it the difference in learning rate, for instance? Or maybe the tasks considered are simply easy enough for any model to work; if so, it could be useful to acknowledge to what extent this is true.

1b) Comparison with IL and CL: On a high level, the comparison with IL and CL algorithms is written as conclusive. I suspect that the failure modes of IL and CL that are described are not due to the algorithms themselves, but rather to the construction of invariance statistics or the choice of negative sample sets (the sets of samples among which variance 1 is requested by VICreg). For instance, if variance (or negative sample set) is taken only across time, the variance object identity is expected to collapse. Similarly, if the variance is taken across the object identity, the variance across time can collapse. So I wonder if the failure of IL and CL is induced by the construction of the variance definition.

(2) Prediction error: When compared to the recording of cortical activity in Figure 7. It is not obvious from the figure which latent space we are talking about mathematically. Is the vector  $z$ ,  $c$  or the prediction error  $e$ ? This is rather important from a neuroscientific point of view, because the prediction error  $e$  is expected to explain the neuronal data. On the other hand,

the prediction error  $e$  is only used in the learning algorithm to define the loss function, but it is not the communication medium between the RNN units  $c$  (or with the encoder  $z$ ).

In the brain, since the measurements are recorded as neural activity, they are communication channels between specific units ( $z$  or  $c$ ). It is probably  $c$  or  $z$  that would already explain the oddball prediction error. I believe that other models, like Forward-forward of Nejad et al., have tried quite hard to address this apparent tension. Whether or not this is resolved by RPL, it thinks it would be beneficial to state the problem and clarify how the algorithm addresses or ignores the issue.

(3) Successor representation without value? I believe the term successor representation is historically relevant in a reinforcement learning (RL) setting and has a precise mathematical definition. Without RL, I feel that learning successor representation is conceptually identical to learning a transition matrix (aka, a primitive world model). I therefore wonder if the pitch for high-level framing of the successor representation is appropriately described or trivial.

(4) Learning in RNN: Learning with recurrent networks appears to be a key in this model presented here (it is in the algorithm name). Yet, this aspect of the model and the literature on biologically plausible learning rules for RNN is not really discussed.

<https://doi.org/10.7554/eLife.110895.1.sa3>

### Reviewer #2 (Public review):

This is a very interesting manuscript, which proposes a novel idea on how cortical networks may learn useful representations of sensory stimuli. The model implementing this idea is thoroughly tested in multiple experimental paradigms. The manuscript is very clearly written. I feel it may have a significant impact on our understanding of cortical circuitry.

<https://doi.org/10.7554/eLife.110895.1.sa2>

### Reviewer #3 (Public review):

Summary:

This paper presents Recurrent Predictive Learning (RPL), a self-supervised model conceptually similar to Joint-Embedding Predictive Architecture (JEPA) models. RPL sequentially observes dynamic scenes to predict subsequent observations. A central claim of the work is that the model's trained representations are simultaneously invariant and equivariant to transformations, such as movement properties that emerge without explicit supervision. These representational qualities are demonstrated through three experiments utilizing two simulated datasets and one naturalistic dataset. Furthermore, the latent embeddings are qualitatively compared with neural data, showing that the model reproduces the successor representation observed in human V1 and the local/global oddball effect in the monkey Prefrontal Cortex.

Strengths:

(1) The paper addresses a fundamental question relevant to both computational neuroscience and machine vision: how the brain learns representations that are simultaneously invariant and equivariant to transformations. The manuscript is well-written, easy to follow, and supported by clear visualizations.

(2) While JEPA-style models have recently gained significant traction in the artificial intelligence community, this paper nicely bridges the gap to neuroscience. By framing these architectures as a theory for visual learning in the brain, the authors provide valuable insights into how predictive frameworks can explain cortical processing.

(3) The qualitative alignment with V1 and PFC data is a particularly strong contribution, as it offers a potential mechanistic explanation for observed neural phenomena through the lens of self-supervised learning.

#### Weaknesses:

(1) The central claim, that both invariance and equivariance emerge spontaneously, requires further scrutiny (see Ghaemi et al., *NeurIPS*, 2025; Garrido et al., *arXiv*, 2024). In particular, the synthetic "moving animal" dataset used in this paper may be too simple to fully support this claim. In latent space prediction, a model must predict both the scene content and the dynamics of movement. Because movement (whether ego-motion or external) is often highly uncertain (or multi-modal), predictive models in naturalistic settings often "collapse" toward learning purely invariant representations, ignoring the hard-to-predict dynamics. In the provided simulations, the movements are extremely predictable. In more complex scenarios, the model would likely prioritize content (invariance) over dynamics (equivariance) unless aided by action-conditioning or explicit factor estimation (Zhang et al., *ICLR*, 2026). The authors' results in Figure 5 using naturalistic video seem to reflect this limitation, given the lower performance on the naturalistic videos compared to the synthetic datasets.

(2) The framing of the RPL model as an entirely new theory of representation learning is slightly overstated. The focus on prediction in representation space rather than input space is the defining characteristic of JEPA and various other Self-Supervised Learning (SSL) models, even sequential prediction. While this paper clarifies the connection between these AI frameworks and cortical circuits, the work would be strengthened by more explicitly positioning RPL within the context of existing JEPA-style models and prior SSL theories of the visual system.

(3) A significant challenge in latent-space SSL is avoiding "representational collapse" (where the model provides a trivial constant output). While the paper alludes to JEPA-like solutions, it lacks a detailed explanation (in both the text and the architectural schematics) of the specific technique used to prevent collapse. Consequently, it is difficult to evaluate the authors' claim of "biological plausibility," as the biological equivalents of common machine learning techniques (such as stop gradient) are not discussed.

(4) Recent work has shown that the capacity (size) of the predictor significantly influences the learned representations in a JEPA-type world model (Garrido et al., 2024). In simpler scenarios, a large enough predictor can allow a model to "memorize" dynamics rather than learning generalized equivariant features. It would be beneficial to see how the ratio of predictor size to encoder size affects the emergence of these features.

#### Methodological Clarifications:

(1) The authors mention a contrastive learning comparison but provide few details. Since contrastive learning is primarily a technique to avoid collapse, it would be a more rigorous baseline if implemented within the same architecture as RPL to isolate the effect of the predictive objective.

(2) In the PFC data comparison (Figure 7f), there appears to be a discrepancy where the local and global conditions show nearly identical results in PFC, while different dynamics in the model. It is unclear if this is a visualization error or a genuine model deviation.

(3) The criteria for selecting specific model variables for comparison with V1 versus PFC are not explicitly defined. Clarification is needed on whether the same latent variables were used for both brain regions or if different layers were selected.

<https://doi.org/10.7554/eLife.110895.1.sa1>

## Author response:

### Public Reviews:

#### Reviewer #1 (Public review):

*The paper describes a biologically plausible version of JEPa using recurrent neural networks called RPL for recurrent predictive learning. Given an embedding  $z_t$  a recurrent neural network processes these inputs with the form:  $c_{t+1} = \text{RNN}(c_t, z_t)$ . Then the predictive network  $f$  is predicting the future inputs with the format:  $\min | |f(c_t) - \text{stop grad}(z_{t+\Delta_t})| |^2$ . I understand that a prediction error is defined as:  $e = z_{t+\Delta_t} - f(c_t)$  to model cortical measurements in the oddball task.*

*The RPL model is also shown to build an internal world model, with "real-world" data like the movement of moving animals or speech signals. The representation is then compared to V1 data and expected prediction error signals in an oddball setting. In a stacked hierarchy of RNN learning with RPL, the higher layers appear to learn high-level latent variables, although gradients are not propagated downward to the lower layers.*

*The paper tackles an open question: Self-supervised learning is thought to be a fundamental principle to explain how computation is structured in the brain. Cortical data suggest qualitatively that prediction error is a core principle of representation learning in the brain, but the field is still looking for a simple yet expressive model that would explain how the cortex learns its representations. RPL contributes in that direction by making a useful link between cortical representation learning in RNN models and the JEPa learning algorithm that was demonstrated to scale to large world model learning from video data by Lecun's group. It is very useful to connect this popular deep learning algorithm to cortical data.*

*The model formalism is relatively elegant and simple: Simple next input prediction objectives are conceptually simple but not necessarily trivial to build at scale. There is a clear benefit in comparison with contrastive or IL methods because they are free from dataset-specific data augmentation and negative samples. Thereby moving the comp neuro field towards conceptually simpler models of representation in the cortex. Yet predictive only models (and in particular predictive models in latent space instead of pixel space) are not easy to build in a stable fashion. JEPa family is basically intended to solve this question; it is very nice and timely to bring this to comp neuro.*

*The methodology combining comp neuro and deep learning makes sense: The conceptual and qualitative analogy with cortical prediction errors is relevant and consistent with what is expected as a model of self-supervised learning in cortical models. The methodology to compare RPL with IL and CL is methodologically meaningful and grounded: showing, for instance, how some of the models fail to represent some latent structure in some toy datasets is interesting.*

*(1.1) h-RPL: The h-RPL is perhaps the most creative departure from the JEPa model family. It would be interesting to say more about what was particularly difficult to see in the latent variables emerging in the hierarchical model. I often find it magical that layer-wise learning rules of this type are not learning redundant representations. Any insights why this is not the case here would be potentially insightful.*

We thank the reviewer for this comment. Regarding representational collapse in h-RPL: each local circuit independently applies the same collapse-preventing strategy as the single-level RPL model: namely, the asymmetric prediction architecture combined with the stop-grad operator. Since this mechanism operates locally within each circuit, it is sufficient to prevent collapse at every level of the hierarchy independently (see also our response to Point P1.3).

The more subtle question is why the circuits learn non-redundant rather than identical representations across the hierarchy. We believe two mechanisms are at play here: First, the hierarchical encoder is a stacked convolutional network, meaning that receptive field sizes grow with depth. This architectural inductive bias naturally encourages successive circuits to operate on increasingly spatially integrated features, creating a structural pressure toward learning complementary rather than redundant representations. Second, the growing expressivity of the network with depth means that higher circuits have access to richer, more abstract inputs from which they can extract higher-level latent structure that is not already captured by lower circuits. Together these factors: the local collapse-preventing mechanism and the depth-dependent growth in receptive field size and network expressivity presumably explain why h-RPL builds an increasingly refined and non-redundant representational hierarchy.

What we will do: We will expand our discussion on this point in the revised manuscript. We plan to expand our quantification on how abstractions emerge in h-RPL in future work in which we will also study variations with top-down connections.

*(1.2) In general, I fully support the type of question and ideas that the paper is putting forward. It is, however, very hard in this research field to gain insight into specific conceptual contributions or specific bits of experimental data that the model puts forward. In pointing to the following weaknesses, I am encouraging the authors to lay out more clearly what the unique hypothesis is or the contribution of the RPL model that we should remember it for.*

Thanks for the positive feedback along with the constructive criticism, and we agree that articulating the core contributions more crisply would strengthen the paper.

At its heart, we believe the paper makes two contributions we hope it will be remembered for. First, while prior work has established that invariant representations can be learned via local Hebbianlike learning rules, we show that learning equivariant representations alongside a latent dynamics model requires something qualitatively different: a local circuit; one with recurrent dynamics and an asymmetric predictive architecture. RPL provides a minimal concrete instantiation of this principle.

Second, and perhaps more broadly, the model makes a structural prediction about (cortical) neuronal circuit organization: since the encoder, integrator, and predictor each perform functionally distinct computations, the framework implies the existence of corresponding cell types and connectivity patterns one should look for in experimental data.

What we will do: We will sharpen these above messages in the revised manuscript to ensure these contributions are prominently highlighted throughout the paper.

*(1.3) Comparison with JEPA variants: JEPA variants are integrating different details into the learning algorithm. Integrating, for instance, "masking" of the latent encoder targets, or EMA in the style of BYOL or Siamese networks, for the predicted representations. It is great that RPL does not seem to need any of those (next input prediction is a natural implementation of masking, and EMA does not seem to be used). It is notoriously hard for the JEPA model to work without these features. Since some of these details are sometimes surprisingly crucial for a simulation to work, it would be good to report which of the other important details were key to live without EMA and masking. Is it the*

*difference in learning rate, for instance? Or maybe the tasks considered are simply easy enough for any model to work; if so, it could be useful to acknowledge to what extent this is true.*

We thank the reviewer for raising this important point. There are two key mechanisms that ensure stable, non-trivial training in RPL. First, using a higher learning rate for the predictor relative to the encoder is crucial for stable training. This prevents the predictor from collapsing the encoder representations and was already noted empirically by Chen et al. (2021).

Second, and more fundamentally, predicting at the level of the memoryless encoder output, rather than at the level of the recurrent integrator, is essential to prevent a degenerate solution in which the RNN simply learns to generate an internally predictable time series unrelated to the input. By anchoring the prediction target to the encoder, the model is forced to ground its representations in the sensory input. Intuitively, otherwise the RNN can simply “make up” a predictable time series, which satisfies the learning objective, but would not yield useful internal representations.

Beyond these architectural points, previous work from our group (Srinath Halvagal et al., 2023) has shown mathematically that JEPAs without EMA avoid collapse via an implicit variance regularization mechanism, and we believe RPL benefits from the same principle. Indeed, we now have a more complete theoretical understanding of this, including identifiability proofs for the latent dynamical model under relatively mild assumptions (Mikulasch et al., 2026). This work has recently been accepted at ICML. Other than that, one has to ensure that representations are not already nearly collapsed at the beginning of training. In this paper, we used normalization layers (batchnorm) in the encoder to ensure this.

Finally like all SSL paradigms the augmentation strength is an important hyperparameter that impacts the quality of learned representations. In the temporal predictive setting, the augmentation strength is fixed by the world itself. The only knob we have to play with is the prediction horizon  $\Delta$ . While we typically focused on next-time-step ( $\Delta = 1$ ) prediction, we saw a clear effect in the case of the speech dataset where  $\Delta = 8$ , but not  $\Delta = 1$ , yielded useful representations for the tasks (Fig. 5b).

What we will do: We will discuss the above points more prominently in the discussion to avoid them being overlooked in the methods. Additionally, we will include a plot on the empirical prediction horizon for the speech dataset in the supplementary material for reference.

*(1.4) Comparison with IL and CL: On a high level, the comparison with IL and CL algorithms is written as conclusive. I suspect that the failure modes of IL and CL that are described are not due to the algorithms themselves, but rather to the construction of invariance statistics or the choice of negative sample sets (the sets of samples among which variance 1 is requested by VICreg). For instance, if variance (or negative sample set) is taken only across time, the variance object identity is expected to collapse. Similarly, if the variance is taken across the object identity, the variance across time can collapse. So I wonder if the failure of IL and CL is induced by the construction of the variance definition.*

We thank the reviewer for this thoughtful point. Both RPL and CL implement an implicit variance regularizer by virtue of being JEPAs (Srinath Halvagal et al., 2023), whereas IL uses an explicit regularizer computed along both the batch and time dimensions to avoid representational and dimensional collapse. The failure modes of IL and CL therefore cannot be entirely attributed to the statistics of the input samples chosen for variance regularization, but are instead primarily determined by the choice of prediction and target representations.

What we will do: We will clarify this in the Methods section of the revised manuscript.

*(1.5) Prediction error: When compared to the recording of cortical activity in Figure 7. It is not obvious from the figure which latent space we are talking about mathematically. Is the vector  $z$ ,  $c$  or the prediction error  $e$ ? This is rather important from a neuroscientific point of view, because the prediction error  $e$  is expected to explain the neuronal data. On the other hand, the prediction error  $e$  is only used in the learning algorithm to define the loss function, but it is not the communication medium between the RNN units  $c$  (or with the encoder  $z$ ).*

*In the brain, since the measurements are recorded as neural activity, they are communication channels between specific units ( $z$  or  $c$ ). It is probably  $c$  or  $z$  that would already explain the oddball prediction error. I believe that other models, like Forward-forward of Nejad et al., have tried quite hard to address this apparent tension. Whether or not this is resolved by RPL, it thinks it would be beneficial to state the problem and clarify how the algorithm addresses or ignores the issue.*

Thanks for pointing out the issue with regards to clarity and for raising the important but subtle point about prediction error representation. To answer the immediate question asking which vector we use in Figure 7, it is the vector  $c$  corresponding to the integrator representations. We agree this should be stated explicitly and will update the manuscript accordingly.

On the more general point, we agree that the tension between recordable neural activity and the computational role of prediction errors is an important issue. We do already briefly engage with it in the Discussion (subsection “Relation to previous modeling work”), where we note that under RPL “inter-areal communication is dominated by representations rather than error signals”. However, we agree that this point should be surfaced more directly.

To elaborate, under classical predictive coding, prediction errors are the inter-areal communication channel and are therefore expected to be directly observable in neural recordings, e.g., as oddball responses. Under RPL, this is not the case:  $e$  is computed locally within a circuit and serves only as a learning signal for synaptic plasticity, not as a signal propagated between circuits or areas. What cortex primarily encodes and communicates in our framework are predictive representations, not reconstruction errors. Accordingly, what should map onto recorded population activity are the representations  $c$  (and  $z$ ), while locally computed prediction errors could in principle remain observable as more circumscribed or transient mismatch-like signals within a circuit.

We would like to push this point further. The reviewer frames this as a tension that RPL needs to resolve, but growing neurophysiological evidence suggests that classical residual-difference prediction errors may not be a dominant mode of cortical encoding in the first place. Furutachi, Franklin, et al. (2024) showed that V1 responses to unexpected visual stimuli do not encode how input deviates from predictions, but instead selectively amplify the representation of the unexpected stimulus itself. Very recently, Furutachi and Hofer (2026) generalize this into a revised framework in which feedforward pathways transmit sensory representations modulated by prediction-error magnitude, rather than residual differences. Vasilevskaya et al. (2026) constrain the space of plausible cortical algorithms via functional influence experiments, also concluding that no variant of standard predictive processing is consistent with the full pattern of layer 2/3  $\leftrightarrow$  layer 5 interactions; they propose a JEPA-based model, citing RPL as a promising candidate. The model by Nejad et al. (2025) similarly shares with RPL the property that representations, rather than residual errors, propagate between circuit elements.

Taken together, the apparent tension may be less a problem RPL needs to resolve than one it is well positioned to explain, remaining consistent with the emerging picture of cortex as

encoding amplified sensory features rather than transmitting residual errors across areas.

What we will do: We will add missing information to the main text and sharpen the Discussion with these arguments.

*(1.6) Successor representation without value? I believe the term successor representation is historically relevant in a reinforcement learning (RL) setting and has a precise mathematical definition. Without RL, I feel that learning successor representation is conceptually identical to learning a transition matrix (aka, a primitive world model). I therefore wonder if the pitch for high-level framing of the successor representation is appropriately described or trivial.*

The reviewer makes a valid point on the concept of successor representations. To answer the immediate question, it is not entirely trivial, as we not only observe the emergence of the transition structure (Fig. 6c), but also the encoding of decaying future (but not past) state occupancy (Fig 6d,e). We largely adapted the terminology “successor-like representations” from the study by (Ekman et al., 2023), but we will elaborate a bit further for why we stuck to it. As nicely pointed out by the reviewer, the term “successor representations” was introduced in the RL literature (Dayan, 1993), but further adopted in neuroscience to describe the idea that a neuronal population encodes a predictive representation that reflects the expected future occupancy of future states under a given policy. Ekman et al. (2023) use the term “successor-like representations” to explain the phenomena where the neural activity in V1 (and hippocampus) represent both current and (discounted) future, but not past, state occupancies in a sequence learning task with no explicitly defined policy or value training. In other words, successor-like representations are simply predictive representations.

What we will do: To deal with this dichotomy, we will replace “successor-like representations” with the term “predictive representations” in the abstract and clarify this distinction in the Results section of the revised manuscript.

*(1.7) Learning in RNN: Learning with recurrent networks appears to be a key in this model presented here (it is in the algorithm name). Yet, this aspect of the model and the literature on biologically plausible learning rules for RNN is not really discussed.*

We thank the reviewer for raising this concern. While h-RPL is one step toward more biologically plausible and spatially local learning rules, exploring it further in terms of temporal credit assignment is beyond the scope of the present study and would require a more systematic and in-depth analysis. However, moving toward more biologically plausible learning rules is an interesting research direction that we plan to explore, as we also mentioned in the Discussion (“Limitations and future research directions”).

We think a viable strategy could be to combine a slim spatial credit assignment strategy such as feedback alignment (Nøkland, 2016; Lillicrap et al., 2016) with an online learning rule using eligibility traces for temporal credit assignment such as SuperSpike (Zenke et al., 2018) or e-prop (Bellec et al., 2020). Similar strategies have given promising results for CLAPP (Illing et al., 2021; Zihan et al., 2026).

What we will do: Following the suggestion, we will discuss biologically plausible learning rules for RNNs in the Discussion.

**Reviewer #2 (Public review):**

*This is a very interesting manuscript, which proposes a novel idea on how cortical networks may learn useful representations of sensory stimuli. The model implementing this idea is thoroughly tested in multiple experimental paradigms. The manuscript is very clearly written. I feel it may have a significant impact on our understanding of cortical circuitry.*

**Reviewer #3 (Public review):**

*This paper presents Recurrent Predictive Learning (RPL), a self-supervised model conceptually similar to Joint-Embedding Predictive Architecture (JEPA) models. RPL sequentially observes dynamic scenes to predict subsequent observations. A central claim of the work is that the model's trained representations are simultaneously invariant and equivariant to transformations, such as movement properties that emerge without explicit supervision. These representational qualities are demonstrated through three experiments utilizing two simulated datasets and one naturalistic dataset. Furthermore, the latent embeddings are qualitatively compared with neural data, showing that the model reproduces the successor representation observed in human V1 and the local/global oddball effect in the monkey Prefrontal Cortex.*

*The paper addresses a fundamental question relevant to both computational neuroscience and machine vision: how the brain learns representations that are simultaneously invariant and equivariant to transformations. The manuscript is well-written, easy to follow, and supported by clear visualizations.*

*While JEPA-style models have recently gained significant traction in the artificial intelligence community, this paper nicely bridges the gap to neuroscience. By framing these architectures as a theory for visual learning in the brain, the authors provide valuable insights into how predictive frameworks can explain cortical processing.*

*The qualitative alignment with V1 and PFC data is a particularly strong contribution, as it offers a potential mechanistic explanation for observed neural phenomena through the lens of self-supervised learning.*

*(3.1) The central claim, that both invariance and equivariance emerge spontaneously, requires further scrutiny (see Ghaemi et al., NeurIPS, 2025; Garrido et al., arXiv, 2024). In particular, the synthetic "moving animal" dataset used in this paper may be too simple to fully support this claim. In latent space prediction, a model must predict both the scene content and the dynamics of movement. Because movement (whether ego-motion or external) is often highly uncertain (or multi-modal), predictive models in naturalistic settings often "collapse" toward learning purely invariant representations, ignoring the hard-to-predict dynamics. In the provided simulations, the movements are extremely predictable. In more complex scenarios, the model would likely prioritize content (invariance) over dynamics (equivariance) unless aided by action-conditioning or explicit factor estimation (Zhang et al., ICLR, 2026). The authors' results in Figure 5 using naturalistic video seem to reflect this limitation, given the lower performance on the naturalistic videos compared to the synthetic datasets.*

We thank the reviewer for the feedback. We agree that further validation on more complex datasets would strengthen the claims, and we take this point seriously. If the reviewer has any suggestions for a specific alternative dataset, we would welcome any recommendations.

Regarding the mouse video data specifically, we realized that this is a suboptimal benchmark rather than a shortcoming of our method. The culprit presumably is that the mice remain largely stationary, leading to a heavily imbalanced velocity distribution peaked near zero (Supplementary Fig. S9 [↗](#)). This imbalance makes equivariance evaluation unreliable regardless of the learning algorithm. For example, end-to-end supervised training results in an  $R^2$  of 0.19 compared to  $0.08 \pm 0.02$  for RPL.

Regarding the moving animal dataset, we note that the dynamics are not trivial from an SSL perspective: unlike moving MNIST (Srivastava et al., 2015), the dataset includes changes in scale and orientation, both features that invariance-focused SSL models can easily ignore, yet RPL recovers reliably. For example, this discrepancy can be seen in [Supplementary Table](#)

S1 [↗](#) where we compare to InfoNCE and CPC. That said, we acknowledge the reviewer's broader concern and will seek to validate RPL on more complex datasets.

While it would be nice to compare to related work by Ghaemi et al. (2024), this study used 3DIEBench (Garrido et al., 2023). Unfortunately, 3DIEBench's reliance on pair-based representations with annotated but random augmentations (such as rotations or color changes) precludes the possibility of smooth latent traversals that would be required for RPL to learn from the same dataset. We will look into whether it is computationally feasible to adapt or regenerate a similar dataset that meets the requirements for temporal prediction.

Regarding stochasticity, we agree that predictive learning in latent space is most natural in approximately deterministic settings, whereas real world sensory information often comprises non-deterministic elements. While a deeper treatment of such stochastic environments is beyond the scope of the present manuscript, it will be the focus of ongoing and future work. Regarding ongoing work, it is worth mentioning that in recent work from our group (Hauri et al., 2026), we have demonstrated that RPL's core objective can replace the reconstruction loss in Dreamer, achieving competitive performance in complex, stochastic environments. While we did not systematically evaluate equivariance in this study, the results suggests that representation-space predictive learning is viable beyond the deterministic regime.

What we will do: We will make the point about the real-world mouse video dataset being a poor benchmark and include the additional  $R^2$  values to show that. Further, we will try to identify or generate alternative datasets to back the equivariance claims and discuss our findings in the light of previous work, e.g., Ghaemi et al. (2024). Moreover, we will sharpen our discussion of our model's limitations in stochastic settings and highlight notable connections to related work.

*(3.2) The framing of the RPL model as an entirely new theory of representation learning is slightly overstated. The focus on prediction in representation space rather than input space is the defining characteristic of JEPA and various other Self-Supervised Learning (SSL) models, even sequential prediction. While this paper clarifies the connection between these AI frameworks and cortical circuits, the work would be strengthened by more explicitly positioning RPL within the context of existing JEPA-style models and prior SSL theories of the visual system.*

Thanks for raising this point. We are unsure what the reviewer refers to. We did not frame our work as "an entirely new theory of representation learning," as the reviewer suggests. In fact, we highlight quite the opposite already in the title of our article, which reads: "Understanding neural circuit principles for representation learning through joint-embedding predictive architectures." We do not claim novelty over JEPA as an ML paradigm, we adopt it precisely because it provides a principled, non-generative framework for predictive representation learning, and our goal is to develop a circuit level instantiation that accounts for neural circuit computation. We already discuss a body of previous work of self-supervised learning and JEPAs at length. Since the reviewer did not specify what they are missing, we will briefly reiterate what is already there.

Our contribution is a theory of representation learning in the brain, built on JEPAs as the underlying ML framework. The Title and Introduction already position our work quite explicitly this way. Specifically, we mention prior work on JEPAs (CPC, BYOL, SimSiam, I-JEPA, seq-JEPA, V-JEPA, V-JEPA 2), while noting that "most JEPAs developed in machine learning are poor models of cortical computation" because of their reliance on negative sampling, transformers, masking, static images, and/or known parametrized transformations, and motivate RPL as the minimal candidate that "must instead rely on recurrent neural dynamics, learn from streaming sensory input without masking, support both invariant and equivariant representations, and reproduce key neurophysiological observations."

The Discussion (“Relation to previous modeling work”) further details the specific novelties of RPL relative to existing sequential JEPA-style and SSL models like CPC (Oord et al., 2018), V-JEPA (Bardes et al., 2024), V-JEPA 2 (Assran et al., 2025), seq-JEPA (Ghaemi et al., 2024). In brief:

RPL is a recurrent JEPA based on RNN dynamics, not transformers, and learns from streaming sensory input without masking or random negative sampling;

It explicitly compares three prediction-error topologies (RPL vs. invariance learning vs. contextprediction; Fig. 2, Suppl. Fig. S2, S6) and shows that asymmetric recurrent prediction is essential for jointly learning invariant and equivariant representations;

Importantly, it does so via pure temporal prediction without access to underlying transformations, a property shared by very few JEPAs. The closest exception is VJ-VCR (Drozdov et al., 2024) which uses an explicit variance-covariance regularization (VCRReg) in a JEPA, which we will cite in the revised manuscript;

It provides the first hierarchical JEPA optimizing local prediction errors at multiple levels (h-RPL, Fig. 8), as envisioned by LeCun (2022) but not previously implemented;

It connects directly to neurophysiological data: successor-like representations in human V1 and abstract sequence representations in macaque PFC, which provides qualitative correspondence between JEPA components and cortical activity that the existing JEPA literature, focused on ML benchmarks, does not address.

Finally, our article already includes a discussion paragraph on recent self-supervised learning models in the context of the brain where we discuss work by Nejad et al. (2025) and Asabuki et al. (2025). Most other SSL theories of the visual system rely on static images and recognition tasks (Yerxa et al., 2024; Margalit et al., 2024). However, there are two studies that include temporal prediction objectives and are worth mentioning with more details: First, Bakhtiari et al. (2021) show that representations similar to ventral and dorsal pathways in the visual system can emerge in a two-pathway encoder architecture within the CPC model. Second, Niu et al. (2024) use a “straightening” objective together with VCRReg as a practical model of the perceptual straightening hypothesis (Hénaff et al., 2019). Though not a JEPA (i.e., has no predictor network), it can decode equivariant factors in a sequential MNIST dataset where only single factors change throughout a video.

What we will do: We will carefully review our discussion of previous work and further discuss Drozdov et al. (2024), Bakhtiari et al. (2021), and Niu et al. (2024) in the revised manuscript.

*(3.3) A significant challenge in latent-space SSL is avoiding “representational collapse” (where the model provides a trivial constant output). While the paper alludes to JEPAlite solutions, it lacks a detailed explanation (in both the text and the architectural schematics) of the specific technique used to prevent collapse. Consequently, it is difficult to evaluate the authors’ claim of “biological plausibility,” as the biological equivalents of common machine learning techniques (such as stop gradient) are not discussed.*

Thanks for pointing this out. Our model avoids collapse through the asymmetric stop-grad / predictor architecture. It does not require an EMA, when the predictor learns with a faster learning rate than the rest of the network (see also our response to Point P1.3).

The use of stop-grad suggests that a circuit learning with RPL needs to compute a vector-based instructive learning signal. While we do not explicitly model the circuit level mechanisms of how this could be implemented in the brain, excitation-inhibition balance is one possibility (Rossbroich et al., 2025). Finally, differences in learning rate can be implemented both structurally or functionally in the brain (see Liu et al. (2025) for instance),

or activity normalization is suggested as a canonical computation in biological neural circuits (Carandini et al., 2012).

What we will do: We will make sure to discuss these putative biological mechanisms in the revised manuscript.

*(3.4) Recent work has shown that the capacity (size) of the predictor significantly influences the learned representations in a JEPA-type world model (Gorrindo et al., 2024). In simpler scenarios, a large enough predictor can allow a model to “memorize” dynamics rather than learning generalized equivariant features. It would be beneficial to see how the ratio of predictor size to encoder size affects the emergence of these features.*

Thanks for raising this concern. We don't observe noticeable difference in position and velocity decoding when changing the width or depth of the MLP predictor in the moving animals data. However, performance on rotation speed and orientation decoding scales with the changes in width, but not depth of the predictor. This analysis excludes the effect of integrator's capacity as it directly affects the dimensionality of the representations, even though it also effectively contributes to prediction computation in RPL.

What we will do: We will include a figure how how task performance varies with the predictor's width and depth.

#### *Methodological Clarifications*

*(3.5) The authors mention a contrastive learning comparison but provide few details. Since contrastive learning is primarily a technique to avoid collapse, it would be a more rigorous baseline if implemented within the same architecture as RPL to isolate the effect of the predictive objective.*

Thanks for the question. We already use the same network model as in RPL for the contrastive predictive learning (InfoNCE) baseline in [Supplementary Table S1](#) and mentioned in the main text (l.164).

What we will do: We will mention the architecture of the non-linear predictor used for InfoNCE baseline in Methods more explicitly.

*(3.6) In the PFC data comparison (Figure 7f), there appears to be a discrepancy where the local and global conditions show nearly identical results in PFC, while different dynamics in the model. It is unclear if this is a visualization error or a genuine model deviation.*

Thanks for picking up on this subtlety in the experimental results. To clarify, it is a model deviation but an interesting one. The local and global responses do look quite similar in the original PFC data. They differ in that the global oddball (xY | xx and xx | xY) response has a secondary peak that encodes the presence of the global oddball, whereas the initial response is actually dominated by local oddball encoding (xY vs xx). Concretely, this results in the response to the xx | xY condition only showing up weakly in the data and at a time lag with respect to the initial local oddball response. Our model, however, does not show the transient initial response to local oddballs in the decoding direction for global oddballs. In a sense, the network model encodes the global oddball concept more robustly than is seen in the PFC data. That said, whether this indicates a genuine difference in representational strategies that needs to be further accounted for, or whether it is an issue stemming from limited sub-sampling of PFC neurons, remains unclear.

*(3.7) The criteria for selecting specific model variables for comparison with V1 versus PFC are not explicitly defined. Clarification is needed on whether the same latent variables were used for both brain regions or if different layers were selected.*

To clarify, the successor-like representations in human V1 and abstract representations in macaque PFC are two different experiments, so each has different latent variables requiring different RPL models. The architecture used for each experiment is detailed in Methods and the criteria for selecting each architecture was the simplest that should work given the task complexity. Throughout the paper, all representation analysis is done on the output of integrator (c) unless said otherwise. We hope this resolves the confusion.

## References

- Chen, Xinlei et al. (2021). “Exploring simple siamese representation learning”. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 15750–15758.
- Srinath Halvagal, Manu et al. (2023). “Implicit variance regularization in non-contrastive SSL”. In: Advances in Neural Information Processing Systems 36, pp. 63409–63436.
- Mikulasch, Fabian A et al. (2026). Understanding Self-Supervised Learning via Latent Distribution Matching. arXiv: 2605.03517[cs.LG].
- Furutachi, Shohei, Alexis D. Franklin, et al. (Sept. 2024). “Cooperative thalamocortical circuit mechanism for sensory prediction errors”. en. In: Nature 633.8029. Publisher: Nature Publishing Group, pp. 398–406. doi: 10.1038/s41586-024-07851-w.
- Furutachi, Shohei and Sonja B Hofer (2026). “Rethinking Predictive Processing”. In: Annual Review of Neuroscience 49.
- Vasilevskaya, Anna et al. (2026). “A functional influence based circuit motif that constrains the set of plausible algorithms of cortical function”. In: bioRxiv. doi: 10.64898/2026.01.29.702557. pdf.
- Nejad, Kevin Kermani et al. (July 2025). “Self-supervised predictive learning accounts for cortical layer-specificity”. en. In: Nat Commun 16.1, p. 6178. doi: 10.1038/s41467-025-61399-5.
- Ekman, Matthias et al. (Feb. 2023). “Successor-like representation guides the prediction of future events in human visual cortex and hippocampus”. In: eLife 12. Ed. by Morgan Barense et al., e78904. doi: 10.7554/eLife.78904.
- Dayan, Peter (1993). “Improving generalization for temporal difference learning: The successor representation”. In: Neural computation 5.4, pp. 613–624.
- Nøkland, Arild (2016). “Direct feedback alignment provides learning in deep neural networks”. In: Advances in neural information processing systems 29.
- Lillicrap, Timothy P et al. (2016). “Random synaptic feedback weights support error backpropagation for deep learning”. In: Nature communications 7.1, p. 13276.
- Zenke, Friedemann et al. (2018). “Superspike: Supervised learning in multilayer spiking neural networks”. In: Neural computation 30.6, pp. 1514–1541.
- Bellec, Guillaume et al. (2020). “A solution to the learning dilemma for recurrent networks of spiking neurons”. In: Nature communications 11.1, p. 3625.
- Illing, Bernd et al. (2021). “Local plasticity rules can learn deep representations using self-supervised contrastive predictions”. In: Advances in Neural Information Processing Systems 34.
- Zihan, Wu S et al. (2026). “Can Local Learning Match Self-Supervised Backpropagation?” In: arXiv preprint arXiv:2601.21683.

- Srivastava, Nitish et al. (2015). “Unsupervised learning of video representations using lstms”. In: International conference on machine learning. PMLR, pp. 843–852.
- Ghaemi, Hafez et al. (2024). “Seq-JEPA: Autoregressive Predictive Learning of Invariant-Equivariant World Models”. In: NeurIPS 2024 Workshop: Self-Supervised Learning - Theory and Practice.
- Garrido, Quentin et al. (2023). “Self-supervised learning of split invariant equivariant representations”. In: arXiv preprint arXiv:2302.10283.
- Hauri, Michael et al. (2026). “Dreamer-CDP: Improving Reconstruction-free World Models Via Continuous Deterministic Representation Prediction”. In: arXiv preprint arXiv:2603.07083.
- Oord, Aaron van den et al. (July 2018). “Representation Learning with Contrastive Predictive Coding”. In: arXiv:1807.03748 [cs, stat]. arXiv: 1807.03748.
- Bardes, Adrien et al. (2024). V-JEPA: Latent Video Prediction for Visual Representation Learning.
- Assran, Mido et al. (2025). “V-jepa 2: Self-supervised video models enable understanding, prediction and planning”. In: arXiv preprint arXiv:2506.09985.
- Drozdo, Katrina et al. (2024). “Video representation learning with joint-embedding predictive architectures”. In: arXiv preprint arXiv:2412.10925.
- LeCun, Yann (2022). “A Path Towards Autonomous Machine Intelligence Version 0.9.2, 2022-0627”. en. In.
- Asabuki, Toshitake et al. (2025). “Learning predictive signals within a local recurrent circuit”. In: Proceedings of the National Academy of Sciences 122.27, e2414674122. doi: 10.1073/pnas.2414674122.
- Yerxa, Thomas et al. (2024). “Contrastive-equivariant self-supervised learning improves alignment with primate visual area it”. In: Advances in neural information processing systems 37, pp. 96045–96070.
- Margalit, Eshed et al. (2024). “A unifying framework for functional organization in early and higher ventral visual cortex”. In: Neuron 112.14, pp. 2435–2451.
- Bakhtiari, Shahab et al. (2021). “The functional specialization of visual cortex emerges from training parallel pathways with self-supervised predictive learning”. In: Advances in Neural Information Processing Systems. Ed. by M. Ranzato et al. Vol. 34. Curran Associates, Inc., pp. 25164–25178.
- Niu, Julie Xueyan et al. (2024). “Learning predictable and robust neural representations by straightening image sequences”. In: Advances in Neural Information Processing Systems 37, pp. 40316–40335.
- Hénaff, Olivier J et al. (2019). “Perceptual straightening of natural videos”. In: Nature neuroscience 22.6, pp. 984–991.
- Rossbroich, Julian et al. (2025). “Breaking Balance: Encoding local error signals in perturbations of excitation-inhibition balance”. In: bioRxiv, pp. 2025–05.
- Liu, Peng et al. (2025). “Layer-specific changes in sensory cortex across the lifespan in mice and humans”. In: Nature neuroscience 28.9, pp. 1978–1989.
- Carandini, Matteo et al. (2012). “Normalization as a canonical neural computation”. In: Nature reviews neuroscience 13.1, pp. 51–62.

<https://doi.org/10.7554/eLife.110895.1.sa0>