### nature communications



**Article** 

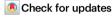
https://doi.org/10.1038/s41467-023-39076-2

# Waves traveling over a map of visual space can ignite short-term predictions of sensory input

Received: 19 August 2022

Accepted: 25 May 2023

Published online: 09 June 2023



Gabriel B. Benigno<sup>1,2,3</sup>, Roberto C. Budzinski<sup>1,2,3</sup>, Zachary W. Davis  $^4$ , John H. Reynolds  $^4$  & Lyle Muller  $^{1,2,3} \boxtimes$ 

Recent analyses have found waves of neural activity traveling across entire visual cortical areas in awake animals. These traveling waves modulate the excitability of local networks and perceptual sensitivity. The general computational role of these spatiotemporal patterns in the visual system, however, remains unclear. Here, we hypothesize that traveling waves endow the visual system with the capacity to predict complex and naturalistic inputs. We present a network model whose connections can be rapidly and efficiently trained to predict individual natural movies. After training, a few input frames from a movie trigger complex wave patterns that drive accurate predictions many frames into the future solely from the network's connections. When the recurrent connections that drive waves are randomly shuffled, both traveling waves and the ability to predict are eliminated. These results suggest traveling waves may play an essential computational role in the visual system by embedding continuous spatiotemporal structures over spatial maps.

Five percent of synapses received by a neuron in the visual cortex arrive through the feedforward (FF) pathway that conveys sensory input from the eyes<sup>1-4</sup>. While these FF synapses are strong<sup>5</sup>, "horizontal" recurrent connections coming from within the cortical region make up about 80% of total synaptic inputs, with 95% of these connections arising from a very local patch (2 mm) around the cell<sup>4</sup>. The anatomy of the visual system thus indicates that cortical neurons interact with other neurons across the retinotopically organized maps<sup>6</sup> that assign nearby points in visual space to nearby points in a cortical region via these horizontal connections. Models of the visual system predominantly focus only on FF<sup>7,8</sup> and feedback (FB)<sup>9</sup> connections. One result of this focus is that, in models of the visual system, neurons in the visual cortex are often modeled as non-interacting "feature detectors" with fixed selectivity to features in visual input (driven by FF connections) that can be modulated by expectations generated in higher visual areas (driven by FB connections). Neuroscientists have long been interested in how horizontal connections shape neuronal selectivity<sup>10,11</sup> and "non-classical" receptive fields<sup>12-16</sup>. More recently, neuroscientists have also been interested in adding these connections to deep learning models to understand neuronal selectivity in the visual cortex<sup>17,18</sup>. It remains unclear, however, how horizontal connections shape the moment-by-moment computations in the cortex while processing visual input.

Recent analyses of large-scale recordings have revealed that horizontal connections profoundly shape spatiotemporal dynamics in the cortex. Traveling waves driven by horizontal connections have been observed in the visual cortex of anesthetized animals<sup>19-24</sup>. The relevance of traveling waves had previously been called into question, as they were thought to disappear in the awake state<sup>25</sup> or to be suppressed by high-contrast visual stimuli<sup>22,26</sup>. Recent analyses of neural activity at the single-trial level, however, have revealed spontaneous<sup>27</sup> and stimulus-evoked<sup>28</sup> activity patterns that travel smoothly across entire cortical regions in awake, behaving primates during normal vision. These neural traveling waves (nTWs) shift the balance of excitation and inhibition as they propagate across the cortex, sparsely modulating spiking activity as they pass<sup>29</sup>. Because they drive

<sup>1</sup>Department of Mathematics, Western University, London, ON, Canada. <sup>2</sup>Brain and Mind Institute, Western University, London, ON, Canada. <sup>3</sup>Western Academy for Advanced Research, Western University, London, ON, Canada. <sup>4</sup>The Salk Institute for Biological Studies, La Jolla, CA, USA.

e-mail: lmuller2@uwo.ca

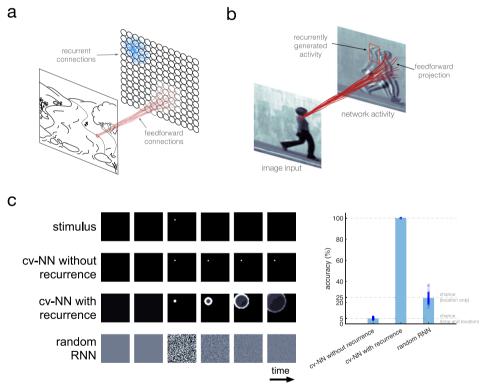
fluctuations in neural excitability<sup>27,30</sup>, nTWs show that neurons at one point in a visual area (representing a small section of visual space) can strongly interact with neurons across the entire cortical region. These results thus indicate that cortical neurons may share information about visual scenes broadly across the retinotopic map through nTWs generated by horizontal connections.

What computations, then, can be done with waves of neural activity traveling across a map of visual space? To address this question, we studied a complex-valued neural network (cv-NN) processing visual inputs ranging from simple stimuli to natural movies. In these networks, activity at each node is described by a complex number. Complex numbers extend the arithmetic of the real number system, and as with standard, real-valued neural networks, nodes receive inputs based on connection weights, with the activity of each node determined by an activation function. The network state is then described by a vector of complex numbers, each element of which can represent the activation of a small patch of neurons in a single region of the visual cortex<sup>31,32</sup>. cv-NNs exhibit similar or superior performance to standard, real-valued neural networks in many supervised learning tasks<sup>33</sup> and have been used effectively in explaining biological neural dynamics<sup>34</sup>. Here, we modified the standard FF architecture used in deep learning and computer vision to include horizontal recurrent connections, where neurons in a single processing layer form a web of interconnections similar to the horizontal connections in the visual

cortex. Horizontal recurrent connections are thought to provide advantages<sup>17</sup> over the standard FF architecture used in computer vision tasks<sup>8,35</sup>; however, current methods for incorporating recurrent horizontal fibers to convolutional network models of the visual system severely limit both the time window over which recurrent activity can be considered and the ease with which the networks can be trained<sup>17</sup>. In recent work, we have introduced a mathematical approach to understand the recurrent dynamics in a specific complex-valued model<sup>36</sup>. Here, we leverage this understanding to train recurrent complexvalued networks to process visual inputs, ranging from simple stimuli to naturalistic movie scenes. The resulting networks can predict learned movies many frames into the future, entirely from their internal dynamics alone, without external input. During prediction, the recurrent network exhibits prominent nTWs, ranging from simple waves propagating out from a small local input28 to complex traveling wave patterns<sup>37</sup>, raising the possibility that nTWs enable continuous predictions of dynamic and naturalistic visual input.

#### Results

The cv-NN consists of an input layer sending movie frames to a recurrently connected neural network. An individual movie frame, serving as input to the network, is represented by a two-dimensional grid of pixels (input frame, Fig. 1a), and each pixel projects to the recurrently connected layer through FF connections (red lines, Fig. 1a).



**Fig. 1** | **A topographic recurrent network model encodes spatiotemporal information of video frames via internal wave activity. a** Schematic of the complex-valued neural network (cv-NN) model. Nodes (circles) are arranged on a two-dimensional grid and are recurrently connected (blue) locally in space like the cortical sheet. A natural image input projects locally into the network via feedforward connections (red), mimicking retinotopy. **b** Example dynamic of the network model. Due to the spatially local projection of the input image, an imprint of the image is visible in the grid of network activity. Due to the local recurrent connectivity, intrinsic wave activity is generated alongside the input projection. **c** Top row: In a sequence of six frames, exactly one of the first five contains a point stimulus, and the other frames do not. These frames are sequentially input to the network. Second row: When the cv-NN has no recurrence, the stimulus projection remains stationary. Third row: With recurrence, from the time of stimulus, cv-NN

activity contains a projection of the stimulus and a wave radiating outward. Fourth row: Activity in a randomly connected recurrent neural network (RNN) following stimulus onset has a spatially disorganized structure, reflecting its lack of topography and distance-dependent time delays. Right: A linear classifier that received the final network state in the no-recurrence case could not predict the time or location beyond chance-level accuracy (5% overall), and in the random-RNN case, could predict the time but not the location beyond chance (25% overall). In contrast, using the classifier with the sixth with-recurrence network state allowed 100% accuracy since the feedforward projection of the point stimulus triggered a radiating wave that encoded the time and location of the stimulus in the subsequent network states. N=100 trials for each group. Mean  $\pm$  standard deviation of 5.09  $\pm$  0.94, 100  $\pm$  0, and 24.42  $\pm$  4.13, respectively. Source data are provided as a Source Data file.

The recurrently connected layer is arranged on a two-dimensional grid, analogous to the retinotopic arrangement of neurons in visual regions. Horizontal interconnections within the cv-NN then drive recurrent interactions in the network (blue lines, Fig. 1a). Both FF and horizontal recurrent projections in the cv-NN are matched to the approximate scale of connectivity in visual cortex<sup>38,39</sup> so that a single pixel in an input movie drives a local patch of neurons, with overlapping horizontal connections, in the cv-NN. Lastly, nodes in the recurrent layer communicate with time delays approximating axonal conduction speeds along horizontal fibers<sup>40</sup>, which have recently been shown to shape spiking neural activity into nTWs<sup>29</sup>. The combination of FF input and dense interconnections generates complex patterns of activity in the recurrent layer (Fig. 1b). Here, we focus on these recurrent activity patterns to understand their computational role for movie inputs ranging from simple to complex.

## nTWs can simultaneously encode stimulus position and time of onset over spatial maps

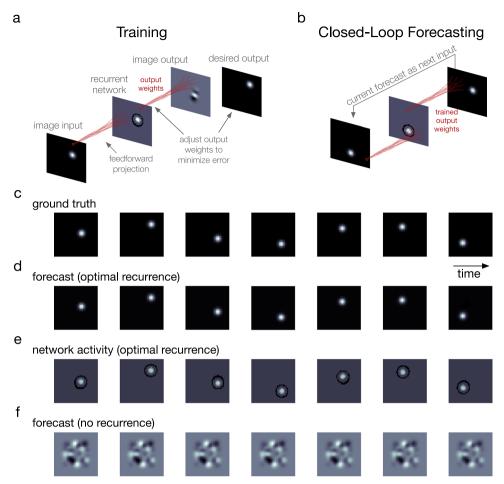
To illustrate how nTWs propagating over sensory maps could facilitate visual computation, we first studied the dynamics generated in response to a single point stimulus. Without recurrent connections, a short point stimulus generates a small bump of activity that remains centered on the point of input ("cv-NN without recurrence", Fig. 1c). With recurrent connections, however, the point stimulus generates a wave that propagates out from the point of input ("cv-NN with recurrence", Fig. 1c). We then studied these stimulus-evoked waves, which are similar in form to those previously observed in the visual cortex of awake primate<sup>28</sup>, in a simple decoding task. Specifically, we let the point stimulus appear at a random time and stimulus location in a series of input frames and then trained a linear classifier to decode the time and location of stimulus onset from the network activity at the final frame. As expected, in the cv-NN without recurrent connections, the classifier performed at chance-level accuracy in this task (Fig. 1c, right; "Methods"-"Stimulus prediction task"). With recurrence, however, the classifier selects the correct time and location of stimulus appearance from the final network state with 100% accuracy. Finally. while standard recurrent neural networks (RNNs) can encode time<sup>41</sup>, an RNN with random connections (and hence lacking the local connectivity and distance-dependent time delays in the cv-NN) also performs at chance level in this task, which requires decoding both stimulus location and onset time (Fig. 1c). This simple illustration shows that traveling waves of neural activity when propagating on an orderly retinotopic map can simultaneously encode stimulus location and onset time, even after the stimulus is no longer present.

#### nTWs aid forecasting movie inputs from simple to complex

Can nTWs enable the processing of the complex, dynamic, and nonstationary visual scenes that we encounter in our natural experience? We approached this question in several steps. We first asked whether, given an input frame from a movie, the cv-NN could be trained to accurately predict the following frame. To perform this more complicated task, we introduced a learning rule that requires training only a linear readout of the recurrent layer (Fig. 2a). This procedure is analogous to a complex-valued implementation of the reservoir computing paradigm<sup>42</sup>, which has recently found wide applications in nonlinear dynamics and physics. In the reservoir computing framework, an input signal drives activity in a recurrently connected layer. Activity in the recurrent layer is then decoded by a set of output weights, which are trained to produce a target output signal. Because of both its efficacy and relative efficiency in training, this framework has proven promising for learning predictive models of chaotic systems<sup>43,44</sup>, and reservoir computing has recently been used to learn and predict a range of important systems in physics<sup>45,46</sup>. This training process, however, has never before been applied to naturalistic movie scenes. We find the cv-NN can be reliably and efficiently trained to predict the next frame in a movie input (Supplementary Table 2, Moving Bump Input). With a cv-NN trained on a movie, the predicted next frame can then be provided as input in place of the original movie (Fig. 2b). Recent work on neural networks for processing movies has focused on predicting the next frame in a video sequence based on training on a large database of inputs<sup>47-49</sup>. In some cases, these predictions can then be fed back as input, allowing the network to recursively generate predictions from its own internal weights<sup>50–58</sup>. We will call this process, where during prediction, a network receives no external movie input and generates future predictions solely from its internal structure, closed-loop forecasting (CLF). Previous work has developed networks that can perform accurate CLF on the order of ten frames into the future<sup>50-58</sup>, with predicted frames becoming increasingly blurry. In this work, we asked a cv-NN to learn and perform CLF on individual movies. We find that cv-NNs trained on an individual movie can self-generate sharp forecasts of that movie many (25-100) frames into the future while receiving no external input. This system can be seen as a simple dynamical autoencoder, where a few input frames can ignite the self-generation of successive frames from its internal dynamics alone. This provides a framework that can give insight into how the visual system could create predictions by continuously changing weights based on its sensory input to make shortterm extrapolations into the near future. The cv-NN is an effective model for closed-loop forecasting of entire visual scenes, generating accurate forecasts for movies of a few thousand pixels per frame using only a few thousand recurrently connected nodes.

The visual cortex readily processes and operates on dynamic visual inputs on timescales of milliseconds to seconds. We then asked whether closed-loop forecasting in this system could work on the scale of tens to hundreds of frames in an input movie. Starting with the first half of a movie containing a simple moving bump stimulus tracing out a trajectory in two-dimensional space (Fig. 2c), we find that the trained cy-NN can produce the entire second half of the movie as output from its trained synaptic weights alone (Fig. 2d and Supplementary Movie 1). As in the previous example, activity in the recurrent layer exhibits a dynamic spatiotemporal pattern extending beyond the immediate FF imprint of the stimulus and structured by the recurrent connections in the network (Fig. 2e and Supplementary Movie 1). These results demonstrate that recurrent cv-NNs can produce simple video inputs from their recurrent connections through this training process. Finally, when we remove the recurrent connections, the cv-NN produces an activity pattern that represents only the average of FF stimulus imprints without having learned the underlying spatiotemporal process<sup>47</sup>. In this case, the cv-NN no longer produces an accurate closed-loop forecast (Fig. 2f). These results demonstrate the importance of both the spatiotemporal patterns in the cv-NN and the horizontal recurrent dynamics generating them.

We find that closed-loop forecast performance in this system depends on two key factors: (1) the ratio of horizontal recurrent strength to feedforward input strength and (2) the spatial extent of the recurrence. To study the first factor in detail, we measured closed-loop forecast performance using an index of structural similarity (SSIM)<sup>59</sup>, which quantifies the perceptual match between two images, ranging between 0 (perfect mismatch) and 1 (perfect match). A threshold on the SSIM, determined through test comparisons between an original and noise-corrupted version of a movie, then provides a quantitative criterion for a successful closed-loop forecast (see Supplementary Fig. 2). We studied SSIM between movie frames produced by the closed-loop forecast process and the ground truth at different ratios of recurrence to input (Fig. 3a; see also Supplementary Fig. 1 and "Methods"—"Network connectivity" and "Network dynamics"). Once the stimulus is removed and the closed-loop forecast begins (video frame 1, Fig. 3a), forecast performance in cv-NNs with low recurrent strength quickly drops close to zero (light blue line, Fig. 3a). By contrast, cv-NNs at optimal recurrent strength sustain closed-loop



**Fig. 2** | The network can forecast a simple video input many frames into the future. **a** As in the classification example (Fig. 1), a video frame projects into the network in a spatially local manner, and a recurrent network interaction occurs, generating internal wave activity on top of the projection. The network outputs an image from its network state via a matrix of trainable weights. Training entails one-shot linear regression between a set of network states and the corresponding desired output frames (the one-step-ahead next frames). Shown: a schematic representation of the one-shot linear regression for one time step. **b** Once training

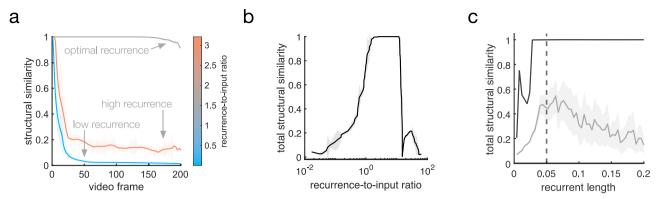
of the readout weights is complete, closed-loop forecasting begins. To properly test how well the network model learned the underlying spatiotemporal process from the training data, it is deprived of ground-truth data of any kind during this step. Instead, the forecast next frame at one time step serves as the input frame for the following time step. **c** Video frames of the data: a bump tracing an orbit. **d** Corresponding closed-loop forecasts generated by the network model with optimal recurrence. **e** Network activity for the optimal-recurrence case. Cosine of phase of activation is shown. **f** Closed-loop forecast in the case without recurrence.

forecasts for long timescales (gray line, Fig. 3a), extending beyond 100 video frames into the future. Importantly, networks where recurrence is too strong also perform poorly, with SSIM dropping near zero within a short timeframe (copper line, Fig. 3a). Systematic quantification of SSIM across ratios of recurrent strength to input strength reveals that performance is best when the recurrence and input are approximately balanced (Fig. 3b), in general agreement with the ratio of feedforward to recurrently generated synaptic drive in visual cortex<sup>60,61</sup>. We next studied performance as a function of the spatial extent of recurrent connectivity. The best performance occurs for recurrent lengths on approximately the same spatial scale as the moving bump stimulus (Fig. 3c), with performance dropping for recurrent lengths outside this range. This result demonstrates that recurrent connections aid closedloop forecasting when matched to the spatial scale of the input. Horizontal recurrent connections in single visual regions span many different retinotopic scales<sup>9,62</sup>, which could enable processing stimuli at multiple spatial scales or moving stimuli with changing scales by the visual system.

The visual system readily processes richly textured and naturalistic visual scenes. To examine this type of stimulus in the cv-NN, we considered naturalistic video inputs for next-frame prediction and closed-loop forecasting. To do this, we used videos from the

Weizmann Human Action Dataset<sup>63</sup>. As above, we trained linear readout weights of the cv-NN on these individual naturalistic movie inputs (Fig. 4a) and then tested whether, given the first half of the input movie, the network could produce the second half in a closed-loop forecast (Fig. 4b). Even with a much more sophisticated input than the previous examples, the cv-NN can be trained rapidly and efficiently on the natural movie inputs (Supplementary Table 2, Walking Person Input). As in previous examples, at optimal values of the network parameters ("Methods"-"Parameter optimization"), the cv-NN accurately produces the natural movie using only its connection weights (Fig. 4c, d and Supplementary Movie 2). In this case, the recurrent connections in the cv-NN create complex wave patterns (Fig. 4e and Supplementary Movie 2). The recurrent connections and their resulting complex activity patterns are important for success in this task, as networks without recurrence do not produce accurate closed-loop forecasts (Fig. 4f).

We then studied what specific features of the recurrent connections enable predicting naturalistic movie inputs. As in the moving bump example, networks perform best when recurrence and input are approximately balanced, and the performance quickly decays when the recurrence is too weak or too strong (Fig. 5a, b). This result shows that, as in the simple case of the moving bump, the complex



**Fig. 3** | **Moving bump forecast performance depends on specific properties of the recurrent connections. a** Structural similarity (SSIM) between a forecast frame and the ground truth as a function of the closed-loop forecast video frame. Each curve corresponds to a different network parameter implementation. Curves have been smoothed by a moving-average filter (filter width of 30 time steps). Shaded error is the absolute difference between filtered and unfiltered. **b** Total structural similarity, in which a single SSIM is calculated for the whole movie as a function of the recurrence-to-input ratio. In the parameter space, each point differs only in recurrent strength. Smoothing and error shading is the same as in (a). **c** Total

structural similarity as a function of recurrent length, which is the fraction of the network's side length spanned by one standard deviation of the Gaussian connectivity kernel. In the three-dimensional parameter space comprising the recurrent strength (rs), recurrent length (rl), and input strength (is), averages (n = 89) across rs-is planes at fixed rl were computed (gray curve). Solid gray line: average. The peak coincides with the standard-deviation width of the Gaussian bump stimulus (dashed vertical line). Shaded area: variance. Solid black curve: maximum structural similarity at each recurrent length. Source data are provided as a Source Data file

spatiotemporal predictions generated by the network depend on a sophisticated interplay between input and recurrent connections. We next studied the role of connection topography and distancedependent time delays. To do this, we started with networks that achieve accurate predictions and randomly shuffled both the connections and time delays, a control that removes the two key factors for generating nTWs in large-scale spiking network models<sup>29</sup> that match waves observed in the visual cortex (Fig. 6a). We then compared the closed-loop forecast performance and network activity in the topographic and shuffled cases. In the topographic case, the cv-NN produces accurate predictions and complex traveling wave patterns. as before (Fig. 6b, c). The shuffled versions of the cv-NN, however. produce spatiotemporally unstructured activity in the recurrent layer (Fig. 6d) and do not achieve accurate closed-loop forecasts, even after the cv-NN was retrained (Fig. 6e; see also Supplementary Table 3 and Supplementary Movie 3). This result demonstrates that with all other architectural features of the network held constant, a randomly connected cv-NN that does not produce nTWs cannot be trained to perform CLF using the same procedure that was previously successful. Shuffling only time delays in the cv-NN and then retraining also substantially drops closed-loop forecast performance (decreasing total structural similarity from 0.99 to 0.02). Further, reducing the conduction speed in half and then retraining also results in a substantial drop in performance (from 0.99 to 0.08). These two control analyses demonstrate that successful closed-loop forecasts depend on a range of time delays in the cv-NN. Finally, the specific spatiotemporal structure of the input movie is also important: a cv-NN at the optimal hyperparameters for a natural movie cannot be retrained to do closedloop forecasting on a randomized (phase-shuffled) version of the same movie (Supplementary Table 1), demonstrating that the cv-NN utilizes the specific spatiotemporal correlations in the movie to generate its forecast. Taken together, these results demonstrate that the complex spatiotemporal patterns generated by horizontal recurrent connections in the cv-NN enable performance on next-frame prediction and closed-loop forecasting tasks for sophisticated natural movie inputs.

## The nTW network model is capable of forecasting multiple movies without retraining

We lastly sought to understand whether the cv-NN could perform closed-loop forecasts on multiple movies it had previously learned and switch flexibly with changing inputs. To do this, we implemented a

simple competitive process ("Methods"—"Movie switching") so that the network could adapt its output based on the similarity of its prediction to its input (Fig. 7a). Specifically, output weights for the cv-NN were trained on individual movies ( $\mathbf{V_1}$  and  $\mathbf{V_2}$ , cf. "Training" in Fig. 7a) and stored in an aggregate matrix ( $\mathcal{V}$ , cf. "Switching" in Fig. 7a). When performing a closed-loop forecast, this extended network model can receive new input from this previously learned set, and then rapidly switch to closed-loop forecasting this new movie input within a few frames without any retraining of weights in the individual output matrices  $\mathbf{V_i}$  (Fig. 7b and Supplementary Movie 4). This result demonstrates that the process of closed-loop forecasting, mediated by horizontal recurrent fibers in the network, can generalize to realistic visual conditions with multiple, changing input streams.

#### **Discussion**

In this work, we have introduced a model to understand whether traveling waves generated by horizontal connections in the visual cortex may play a computational role in processing natural visual inputs. By adapting a recurrent neural network model using a specific dynamical update rule and learning rule, this model learns to forecast video inputs ranging from simple visual stimuli to complex natural scenes. We report here a network model that can be trained to produce quantitatively verified closed-loop forecasts of richly textured naturalistic movies many frames into the future. The cv-NN introduced in this work incorporates the spatial topography and time delays important for shaping activity dynamics in single regions of the visual system<sup>29</sup> and provides a potential computational role for waves of neural activity traveling over maps of visual space. Whether similar principles of spatial topography could benefit RNNs, in general, remains open but represents an interesting potential direction for future work. Further, because the recurrent dynamics in the cv-NN are tractable to detailed mathematical analysis<sup>36</sup>, this recurrent network model opens new possibilities for understanding the mechanisms underlying successful predictions studied here and for designing new applications in future work.

Closed-loop forecasting in the cv-NN demonstrates a form of short-term prediction by nTWs that may be relevant to the online processing of continuous sensory input by the visual system. Consider, for example, a batter in the game of baseball facing a pitcher who has just pitched a curveball, now hurtling toward the batter at over 100 miles per hour. In major league baseball, a pitch takes around 400

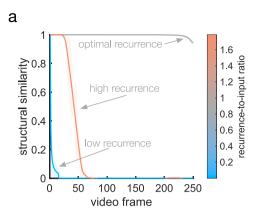


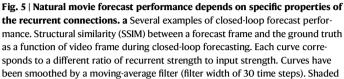
**Fig. 4** | **The recurrent network performs next-frame forecasting of a natural video input. a** Training follows as in the moving bump example (Fig. 2a). **b** Next-frame closed-loop forecasting follows as in the moving bump example (Fig. 2b). **c** Video frames of the data: a person walking, **d** Corresponding closed-loop

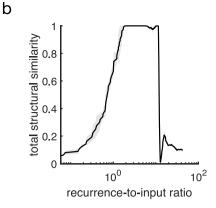
forecasts generated by the network model in the case of optimal recurrence. **e** Corresponding network states for the optimal-recurrence case (**d**). Cosine of phase is shown. **f** Same as (**d**), but in the absence of recurrence.

milliseconds to travel 60 feet from the pitcher's hand to the batter at home plate. Time is required for the neural computations that enable the batter to perceive the ball and estimate its trajectory. This includes both the time required for sensory information to travel from the retinae to relevant brain areas and the time required for computation of the ball's trajectory in space based on these signals. Assuming the entire computation can be accomplished in 150 milliseconds<sup>64</sup>, during this time, the ball will have traveled more than 22 feet. To estimate the likely current location of the ball based on information that was available to the visual system 150 milliseconds ago, the brain may form an internal model of the ball's trajectory in space, informed by previous experience. Consistent with this idea, batters often report that, as the spinning ball travels from the pitcher's mound to home plate, the curveball suddenly changes direction, an illusory percept referred

to as the curveball's "break"<sup>65</sup>. Short-term predictions by nTWs may represent one mechanism for rapid estimation of trajectories, as continuous spatiotemporal structures propagating over the retinotopic map. In this way, closed-loop forecasts in the cv-NN could enable the visual system to estimate the likely trajectory of the ball based on training from the previous visual experience. The curveball's "break" further recalls the process of switching predictions when the input becomes sufficiently discrepant with incoming sensory data (Fig. 7b). When the movie switches from one input to another (top row, "ground truth"), the network generates a transiently indeterminate activity pattern before jumping to the correct forecast (bottom row, "closed-loop forecast"). In this way, the cv-NN may provide a mechanistic framework for specific hypotheses in future work about the interaction of short-term predictions generated by recurrent horizontal fibers and







error is the absolute difference between filtered and unfiltered. **b** Total structural similarity, in which a single SSIM is computed for the whole movie as a function of the recurrence-to-input ratio. In the parameter space, each point differs only in recurrent strength. Smoothing and error shading is the same as in (a). Source data are provided as a Source Data file.

continuously incoming sensory input. The cv-NN could also be useful as a model to explain how the brain encodes, stores, and recovers episodic memories of richly textured visual scenes, which studies of visual search<sup>66</sup> and vivid recollection<sup>67,68</sup> have shown are associated with activity in visual regions.

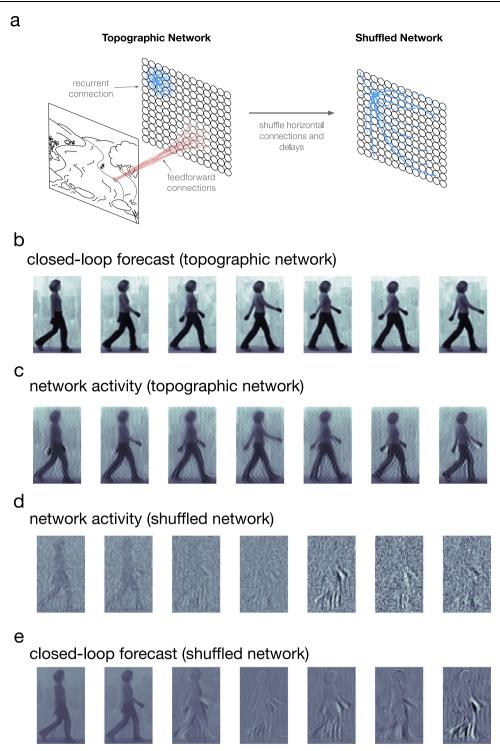
Further, while the cv-NN is not intended to be a veridical simulation of the millions of neurons contributing to nTW dynamics in the visual cortex, this network model is broadly consistent with spatiotemporal dynamics recently observed in the visual system of the alert primate. In the case of a single point stimulus (Fig. 1c), the network produces a traveling wave radiating out from the point of input. This is similar to nTWs detected in single trials during voltage-sensitive dye optical imaging in the primary visual cortex (V1) of awake macagues<sup>28</sup>. nTWs evoked by small visual stimuli (Gaussian spot, with a standard deviation of 0.5° of visual angle) presented during fixation consistently evoked nTWs that propagate over 7.5 mm of V1, representing a significant portion of this cortical area<sup>69</sup>. The spatial extent of the nTWs observed in the experiment provides a point of comparison with the model, as spatial extent determines the scale at which local populations in V1 may influence others across the retinotopic map. In the cv-NN, waves generated by small point stimuli propagate over slightly more than one-third of the network (decaying to half-amplitude after traveling over 37.5% of the network; Fig. 1c). These results demonstrate that nTWs may propagate over broadly similar spatial extents in visual cortex and in the cv-NN.

Another point of comparison with measured neural dynamics centers on the patterns evoked by moving stimuli. In the case of a moving bump stimulus (Fig. 2), the network produces a bump of activity, reflecting FF input driven by the movie but also reflecting recurrently generated activity that extends beyond the feedforward imprint of the stimulus (Fig. 2e). The radius of this recurrently generated activity is approximately twice that of the feedforward bump. This result recalls analyses of Utah array recordings in V1 of awake macaques<sup>70</sup>. Using a moving bar stimulus  $(0.5 \times 4^{\circ})$  of visual angle, moving horizontally at 6.6° per second), the authors found responses in V1 before stimuli entered neurons' classical receptive field (cf. Fig. 2C in ref. 70). The onset times of these anticipatory responses became earlier and earlier along the moving bar's trajectory. These changes in time were confirmed with computational analyses and modeling to be consistent with propagation along horizontal fibers in V1, and the spatial extent of the recurrent interactions is, again, approximately consistent with dynamics during closed-loop forecasting in the cv-NN.

The dynamics of the cv-NN are thus broadly consistent with observations of neuronal dynamics during normal processing in awake, behaving primates. Recent work has demonstrated the importance of the topographic connection patterns and axonal time delays matching those found in the visual cortex to generate nTWs in large-scale spiking network models<sup>29</sup>. Recent theoretical studies have developed complex-valued network models that can provide analytical insight into the time-varying dynamics of spiking neural networks<sup>31,71</sup>, and future work could directly relate dynamics in the cv-NN during movie prediction to the fine-scale spiking dynamics of the networks in the visual cortex. Finally, in the case of naturalistic movie inputs (Fig. 4), the cv-NN produces complex spatiotemporal patterns that can be mathematically described in this model as the summation of multiple traveling waves<sup>36,37</sup>. Future work analyzing large-scale recordings will provide opportunities for comparison between activity patterns in the visual cortex and in the cv-NN during the processing of naturalistic movie inputs.

Another potential extension of the cv-NN is to consider multiple recurrently connected layers with specializations similar to those in different regions of the visual cortex. In this work, we focused on a cv-NN with a single recurrently connected layer to understand the potential computational role of nTWs that have recently been observed in single cortical regions during visual perception in awake animals. nTWs have been observed in many visual areas, including V1<sup>28</sup>,  $V2^{28}$ ,  $V4^{72}$ , and  $MT^{24,27,73}$ . Adding multiple recurrent layers in the cv-NN may provide opportunities in future work for understanding nTW dynamics across visual areas, where spatiotemporal activity patterns have recently been shown to propagate in feedforward and feedback directions in different frequency ranges<sup>74</sup>. Finally, closed-loop forecasts in this cv-NN are not intended to be robust to arbitrary translations or rotations of the visual scene, and adding multiple layers in the cv-NN may provide a degree of translation invariance, which is achieved in CNNs through cascading activity through multiple processing layers<sup>75</sup>, and scale invariance, which may also be made possible through processing in multiple recurrent layers<sup>76</sup>. In this way, extending the cv-NN with multiple recurrent layers represents an important opportunity for understanding the organization and computational role of nTWs occurring in many cortical areas in future work.

These results provide fundamental insight into the function of horizontal recurrent connections, whose effect on the moment-by-moment computations in the visual system has remained unexplained. While there has been much interest in the function of recurrent



**Fig. 6** | **Randomly shuffling recurrent connections eliminates nTWs and the ability to forecast. a** Left: the topographic network model used throughout this study, featuring feedforward projections of the image input (red lines) and local distance-dependent horizontal connectivity (blue lines). There are also synaptic time delays proportional to a node pair's separation distance within the horizontal

recurrent circuitry. Right: by randomizing the horizontal connection weights and time delays, the topography in the network is removed.  $\bf b$  Closed-loop forecasts generated by the topographic network.  $\bf c$  The network activity of the topographic network in response to frames of a natural movie input.  $\bf d$  Network activity of the shuffled network.  $\bf e$  Closed-loop forecasts generated by the shuffled network.

horizontal fibers in the visual cortex, for example, in explaining direction and orientation selectivity in VI<sup>10,11</sup> or in center-surround models of the receptive field<sup>14,16,77</sup>, general computational roles for traveling waves generated by the massive recurrent circuitry in single cortical areas on the single-trial level remain unknown. Successful models of the visual system, including feature-based models and deep convolutional neural networks, have provided insight into how neural

systems could process single image inputs but explain only a fraction of the variance in neural responses to natural sensory stimuli<sup>18,78,79</sup>. Importantly, it is not necessarily the case that all RNNs that can perform CLF will also exhibit nTWs; however, when networks possess the main architectural features found in the visual cortex (local connections, retinotopically ordered inputs, and communication time delays), we have demonstrated that nTWs are tightly linked to CLF. The

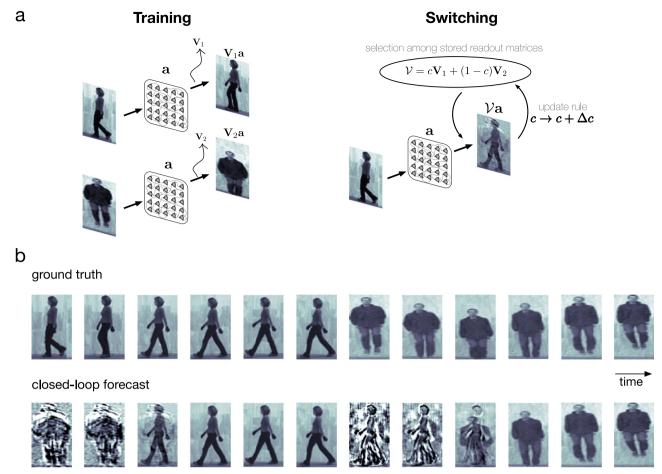


Fig. 7 | The network is capable of forecasting multiple movies without being retrained. a The recurrent network model was adapted to contain a higher-level competitive-learning process. Left: Readout matrices were learned separately for separate examples. Right: Storing the learned readout matrices in an aggregate matrix  $\mathcal{V}$ , the present network state drove the aggregate matrix toward either of the

learned matrices via an unsupervised competitive learning rule. **b** Beginning with feeding frames from movie 1, the network takes some time to recall the learned matrix that results in an accurate closed-loop forecast. Quickly switching to a different movie, the network once again takes some time to adjust its output weights before converging to the correct ones for an accurate closed-loop forecast.

cv-NN may thus provide new opportunities for understanding how the visual system processes continuously updated, movie-like visual inputs, where information is extracted from the visual environment moment-by-moment as it comes from the eye. The sophisticated closed-loop movie forecasts produced by this network, and the fact that this closed-loop forecast process can generalize to multiple movie inputs, represent an important step in explaining the computational role of recurrent connections and traveling waves in the visual cortex.

#### Methods

Custom MATLAB (version R2021a) code was used for all data simulation and analysis in this study.

#### **Network connectivity**

The recurrent network is arranged on a square grid of N nodes. The network grid is treated as a discretized Euclidean plane such that the side lengths span distances of unity. Boundaries are not periodic. The recurrent weight  $w_{ij}$  from node j to node i is inversely proportional to their Euclidean distance  $d_{ij}$  so as to give local connectivity like that of the neocortical sheet. Specifically,  $w_{ij}$  is Gaussian as a function of  $d_{ij}$ :

$$w_{ij} = \alpha \exp\left[-d_{ij}^2/(2\beta^2)\right]. \tag{1}$$

The coefficient  $\alpha$  is called the recurrent strength, and the standard deviation  $\beta$  is called the recurrent length. Both are free parameters. The maximum possible value of  $d_{ii}$  is  $\sqrt{2}$  (corner to corner), and, for

example,  $\beta = 1$  means that the recurrent length equals the network side length. Further, all  $N^2$  such weights are strictly positive, and the N-by-N matrix of such weights is symmetric ( $w_{ij} = w_{ji}$ ). Diagonal weights ( $w_{ii}$ ) are not set to zero.

#### **Network dynamics**

Network dynamics are given by a complex-valued equation. A complex number z is of the form z=x+iy, where x is the real part, y is the imaginary part, and i is the imaginary constant defined as  $i^2 = -1$ . Equivalently,  $z = m \exp[i\phi]$ , where m is the modulus and  $\phi$  is the argument. A complex number is intuitively visualized as a twodimensional vector, where (x,y) is its Cartesian representation and  $(m,\phi)$  is its polar representation. What distinguishes a complex number from a standard two-dimensional vector is the multiplication rule: multiplication of two complex numbers corresponds to both a scaling and a rotation in the so-called complex plane. This property makes complex-valued representations of observable quantities more concise than real-valued representations, and thus, complex numbers are a central tool in physics and engineering. From the perspective of biological vision, a complex-valued representation is useful. Since phase information is important for representing visual inputs, complexvalued models, which efficiently represent phase in the argument  $\phi$ , are ideal. Indeed, complex-valued models of vision are widely explored<sup>80</sup>. Given the practical utility of artificial neural networks and deep learning (including for modeling biological neural networks), complex-valued neural networks, in which the neural activations are

complex-valued, are of great interest. However, they are notoriously difficult to train, especially in a recurrent architecture<sup>32</sup>. We make an advance here on this front by choosing a unique dynamical equation and by exploiting the advantages of reservoir computing.

The discrete-time dynamical equation for each node i is

$$a_i[t+1] = a_i[t] + x_i[t] - i \sum_{j=1}^{N} w_{ij} \exp \left\{ i(a_j[t-\tau_{ij}] - a_i[t]) \right\},$$
 (2)

$$a_i[t+1] := a_i[t+1]/|a_i[t+1]|.$$
 (3)

Here,  $a_i[t]$  is the complex-valued activation at discrete time t,  $x_i[t]$  is the feedforward input of the image stimulus to node i at discrete time t, and  $w_{ij}$  is the recurrent weight from node j to node i ("Methods"—"Network connectivity"). Further,  $\tau_{ij}$  is the discrete time delay between nodes i and j, given by  $\tau_{ij} = \text{round}[d_{ij}/v]$  in which the Euclidean distance  $d_{ij}$  between nodes i and j ("Methods"—"Network connectivity") is scaled by the parameter v, which represents the speed of activation transmission across the network, and  $\text{round}[d_{ij}/v]$  rounds  $d_{ij}/v$  to the nearest integer in accord with the discrete-time dynamics. A v-value of, for example, v = 0.1 means the activation travels a distance of one-tenth the network side length per time step. Lastly, the modulus of  $a_i[t]$  (i.e.,  $|a_i[t]|$ ) is normalized (Eq. (3)), which confines  $a_i[t]$  on the complex unit circle, and thus, the phase of  $a_i[t]$  contains the dynamics. We note that modulus normalization is a common operation used in complex-valued neural networks<sup>32</sup>.

The specific form of Eq. (2) is unique compared to other complex-valued neural-network equations because it involves a pairwise node attraction  $a_j[t-\tau_{ij}]-a_i[t]$ . Another system with pairwise attraction is the Kuramoto model, a popular model for studying synchronization in nonlinear systems<sup>81-83</sup>. Our presented system has a correspondence with the Kuramoto model<sup>84</sup> and allows the description of the dynamics for individual realization in terms of the eigenvalues and eigenvectors of the network<sup>36</sup>. With the described local network connectivity and distance-dependent delays, the presented system gives rise to meaningful spatiotemporal self-organization dynamics.

The initial network state is  $a_i[0] = 0 + 0i$  for all nodes, and the first several time steps contain transient activity associated with the input disrupting the initial steady state of the system. For the stimulus prediction task, this transient activity is important to the model and was used, while for the next-frame forecasting task, it is distracting to the model and was discarded.

#### Image read-in

At each discrete time step, a digital grayscale image is read into the network. Prior to read-in, the image is mean-subtracted and divided by its standard deviation across all its pixels (i.e., z-scored). Image read-in is accomplished with a local feedforward projection, which mimics retinotopy and preserves the spatial correlations in the image. Technically, this is a two-dimensional interpolation using the bilinear kernel common in image processing, which takes a weighted average in the nearest 2-by-2 pixel neighborhood. The projected image has  $\sqrt{N}$  rows and  $\sqrt{N}$  columns like the network grid, and each pixel intensity of the projected image is given by  $x_i[t]$  (Eq. 2). Lastly,  $x_i[t]$  is scaled according to  $x_i[t] := \gamma x_i[t]$ , where  $\gamma$  is called the input strength. In our model,  $\gamma$  is the fourth and final free parameter after the recurrent strength, recurrent length, and conduction speed.

#### Stimulus prediction task

The classification was performed using the basic perceptron. For an input vector  $\mathbf{v} = \begin{bmatrix} 1 v_1 \cdots v_N \end{bmatrix}^T$ , where  $v_1, \ldots, v_N$  are features, and a label  $l \in \{0,1\}$ , the goal is to find a hyperplane  $\mathbf{u}^\mathsf{T}\mathbf{v} = b + u_1v_1 + \cdots + u_Nv_N = 0$ , where  $\mathbf{u} = \begin{bmatrix} b u_1 \cdots u_N \end{bmatrix}^T$  is a vector containing the bias b and weights

 $u_1,\ldots,u_N$ , that separates the data in the N-dimensional feature space according to their binary class (0 or 1). During training, with a suboptimal  $\mathbf{u}$ -vector and one example  $\mathbf{v}$ -vector, the output classification  $l=H(\mathbf{u}^{\mathsf{T}}\mathbf{v})$  is computed, where  $H(\cdot)$  is the Heaviside step function defined as unity for positive argument and zero otherwise. For the desired classification d (either 0 or 1), the signed distance  $\Delta=d-l$  is computed, where  $\Delta\in\{-1,0,1\}$ . With each new example  $\mathbf{v}$ , the  $\mathbf{u}$ -vector is updated using the delta rule  $\mathbf{u}:=\mathbf{u}+\lambda\mathbf{v}\Delta$ , where  $\lambda$  is the learning rate. To use the perceptron in multiclass classification, the one-versus-rest scheme is used. That is, for the set of classes  $C=\{c_1,\ldots,c_M\}$ , binary classification is performed separately M times. Each time i, the two classes are defined such that  $c_i=1$  and  $C\setminus c_i=0$ , where "\" denotes the set difference. Then, there are M weight vectors  $\mathbf{u}_1,\ldots,\mathbf{u}_M$ , and M inner products  $f_1=\mathbf{u}_1^{\mathsf{T}}\mathbf{v}$ , ...,  $f_M=\mathbf{u}_M^{\mathsf{T}}\mathbf{v}$  for a given data vector  $\mathbf{v}$ . The multiclass classification is  $\operatorname{argmax}[f_1,\ldots,f_M]$ .

In the stimulus classification task (Fig. 1c), input frames were 50 by 50 pixels, and the network was 50 by 50 nodes. There were six frames. One of the first five frames was randomly chosen to contain the point stimulus, and the remaining frames were entirely zero intensity. The point stimulus was an isotropic two-dimensional Gaussian of standard deviation of 0.05, and the input frames are defined on the Cartesian grid  $[-2,2]\times[-2,2]$ . The stimulus was centered in one of four equally sized quadrants in the frame. The sequence of frames was sequentially input to the network. There are exactly twenty classes: each of the first five frames times each of the four quadrants in which the point stimulus could occur. The column vector of activations corresponding to the final (sixth) frame was used as predictor for all trials. The task was repeated 100,000 times, with the time of stimulus (1 or 2 or 3 or 4 or 5) and the location of the stimulus (quadrant 1 or 2 or 3 or 4) randomly rechosen each time.

#### **Closed-loop forecasting**

The network outputs an image of  $M_r$  rows and  $M_c$  columns of pixels—the same size as the input image—at each time step. In both examples (moving bump and natural movie), the network was 50 by 50 nodes (N = 2500). Recalling that  $a_i[t]$  is the complex-valued activation of node i at discrete time t (Eqs. (2) and (3)), the output transformation is linear:

$$y_i[t] = \sum_{j=1}^{N} \nu_{ij} a_j[t]'.$$
 (4)

Here,  $y_i[t]$  is the  $i^{\text{th}}$  pixel intensity of the output image, and  $v_{ij}$  is the  $(i,j)^{\text{th}}$  readout weight of the M-by-N matrix  $\mathbf{V}$ , where  $M = M_r M_c$ . The prime notation (') indicates that the activation vector  $\mathbf{a}[t] = \left[a_1[t] \cdots a_N[t]\right]^T$  was mean-subtracted, which was done to avoid an intercept term during training.

The readout weights  $\{v_{ij}\}$  of **V** are the only weights trained in our model, making our network a reservoir computers. Reservoir computers are recurrent neural networks that avoid the issues associated with training recurrent weights and have been shown to perform well in time series forecasting<sup>42</sup>. Suppose training begins at time step 1, after discarding the initial transient, and ends at time step T. Defining  $\mathbf{a}[\mathbf{t}]' = [a_1[\mathbf{t}]' \cdots a_N[\mathbf{t}]']^T$ , the matrix of regressors is then

$$\mathbf{A} = \left[ \mathbf{a}[1]' \dots \mathbf{a}[T]' \right] \tag{5}$$

and the matrix of regressands (desired outputs) is

$$\mathbf{D} = [\mathbf{f}[2] \dots \mathbf{f}[T+1]]. \tag{6}$$

Hence, the desired outputs are simply the set of one-step-ahead frames. Here,  $\mathbf{f}[t]$  is the column vectorization of the  $t^{\text{th}}$  input image frame (before read-in) and is also mean-subtracted. Training entails ordinary least-squares linear regression between  $\mathbf{A}$  and  $\mathbf{D}$ . Because  $\mathbf{D}$  is highly underdetermined (containing far fewer frames than pixels per

frame), the matrix 2-norm of **V** was simultaneously minimized during regression to reduce model bias.

Following training is *closed-loop forecasting*. At this point, the network activation has been primed by being driven with the training frames, and the readout matrix  $\mathbf{V}$  has been trained. In the first time step of closed-loop forecasting, we input the corresponding video frame. Subsequently, for steps  $\{t\}$ , the predicted output at time step t serves as the input for time step t+1.

In the moving bump example (Fig. 2), the frames are 30 by 30 pixels and defined on a  $[-2,2]\times[-2,2]$  Cartesian grid. A two-dimensional isotropic Gaussian of standard deviation 0.2 traced a Lissajous curve given by the parametric equations  $x_c(t) = \sin(t/3)$  and  $y_c(t) = \cos(t/3)$ , where  $(x_c,y_c)$  is the center of the Gaussian in space and t is a continuously valued time variable<sup>85</sup>. The Lissajous trajectory was discretized to have 100 frames per cycle. The first cycle was discarded to omit the initial transient network activity, the network was trained on the subsequent 3 cycles, and closed-loop forecasting was performed on the 2 cycles subsequent to that.

In the natural video example (Fig. 4), a walking video from the Weizmann Human Action Dataset<sup>86</sup> was used, in which a person walks across the scene. We present several key examples here but note that the model successfully performs closed-loop forecasting for all movies in this dataset, where we define a successful closedloop forecast as one in which the total structural similarity is at least 0.9 (Supplementary Fig. 2, Supplementary Table 4). Segmentation masks of the people in the videos are included with this dataset (https://www.wisdom.weizmann.ac.il/~vision/SpaceTimeActions. html). Using these masks, we cropped the frames so that the person was centered throughout the entire walk, giving frames of approximately 80 by 50 pixels. Without performing this step, our network model would fail: the training data would be independent of the closed-loop forecast data since they would occupy exclusive regions of the pixel space, and the model would not generalize to the prediction data. Such nonstationary data have been successfully taught to networks with approximate translation invariance, and translation invariance is likely used in the brain to learn such processes<sup>87</sup>. However, translation invariance is beyond the scope of our study. The frames were then resized to be exactly 80 by 50 pixels. Finally, each video was around 70 frames long. To get more frames without interpolation, we "bookended" each video by concatenating it with its temporal reverse sequence, where one cycle consists of the original frames followed by the bookended frames. The result is a longer video with the same spatiotemporal statistics. The first cycle was discarded to omit the initial transient network activity, the network was trained on the subsequent three cycles, and the closed-loop prediction was performed on the two cycles subsequent to that.

To measure the balance between feedforward input and recurrent interaction, we devised the *recurrence-to-input ratio*. Per Eq. (2), the input and recurrence terms are the column vectors  $\mathbf{x}[t] = [x_1[t] \cdots x_N[t]]^T$  and  $\mathbf{r}[t] = [r_1[t] \cdots r_N[t]]^T$ , respectively, where

$$r_{i}[t] = -i \sum_{i=1}^{N} w_{ij} \exp \left\{ i(a_{j}[t - \tau_{ij}] - a_{i}[t]) \right\}.$$
 (7)

Further, let the matrices

$$\mathbf{R} = \begin{bmatrix} \mathbf{r}[1] \dots \mathbf{r}[T] \end{bmatrix} \tag{3}$$

and

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}[1] \dots \mathbf{x}[T] \end{bmatrix}$$

Table 1 | Intervals over which model parameters were randomly searched during optimization

Parameter	Sampled interval
Recurrent strength	(0, 0.2)
Recurrent length	(0, 0.2)
Input strength	(0, 0.2)
V	(0, 0.1)

be the horizontal concatenations of  $\mathbf{r}[t]$  and  $\mathbf{x}[t]$ , respectively, over closed-loop forecast times  $\{t, t+1, ..., t'\}$ . The ratio is defined as

$$|\mathbf{R}|_{\mathsf{F}}/|\mathbf{X}|_{\mathsf{F}},\tag{10}$$

where  $\|\mathbf{G}\|_F$  denotes the Frobenius matrix norm of a matrix  $\mathbf{G}$ , which is equivalent to the Euclidean vector norm of the vectorization of  $\mathbf{G}$ .

#### **Movie switching**

The network was trained on two movie inputs: one of a walking person (movie 1) and one of a jumping person (movie 1), both from the Weizmann dataset. The same recurrent matrix was used in each case–only the learned matrices ( $\textbf{V}_1$  and  $\textbf{V}_2$  , respectively) differed. Let  $\mathcal V$ =  $c\mathbf{V}_1$  +  $(1-c)\mathbf{V}_2$ , where  $c \in [0,1]$ .  $\boldsymbol{v}$  stores both learned matrices, and the present input modulates the relative contribution of  $V_1$  and  $V_2$  using an update rule for c. The structural similarity between the input and output were computed at each time step t (S[t]), and the change thereof was computed at each time step as  $\Delta S = S[t] - S[t-1]$ . The update rule is  $c := c + \Delta c$ , where  $\Delta c = -\eta \operatorname{sgn}[\Delta S]$  and  $\eta$  is the learning rate, set to 0.1. Depending on which movie (movie 1 or movie 2) drives the network, c tends toward 1 or 0, respectively. Once this happens, this driving input is removed and closed-loop forecasting commences as described. Switching entails instantaneously transitioning from closed-loop forecasting of one movie to driving the network with the frames of another movie. c then updates as described and is followed by closed-loop forecasting again.

#### **Parameter optimization**

The random-search algorithm was used to optimize parameters for closed-loop forecasting. Within specified bounds, each parameter was randomly sampled, giving a point in the parameter space. The parameter space was randomly sampled in this way many times, and each time, the structural similarity index was computed as the performance index. The bounds within which the parameters were sampled are given in Table 1.

#### **Reporting summary**

Further information on research design is available in the Nature Portfolio Reporting Summary linked to this article.

#### Data availability

The point stimulus and moving bump stimulus data generated in this study can be generated from the code available at this study's GitHub repository (https://github.com/mullerlab/benignoEAwavecomp). The raw data of the natural movies used in this study are provided by Lena Gorelick, Moshe Blank, and Eli Shectman of the Weizmann Institute of Science, available at https://www.wisdom.weizmann.ac.il/-vision/SpaceTimeActions.html and this study's GitHub repository. Source data are provided with this paper.

#### **Code availability**

 All codes associated with this study are available at https://github.com/ mullerlab/benignoEAwavecomp<sup>88</sup>.

#### References

- Peters, A. & Payne, B. R. Numerical relationships between geniculocortical afferents and pyramidal cell modules in cat primary visual cortex. Cereb. Cortex 3, 69–78 (1993).
- Latawiec, D., Martin, K. A. & Meskenaite, V. Termination of the geniculocortical projection in the striate cortex of macaque monkey: a quantitative immunoelectron microscopic study. *J. Comp. Neurol.* 419, 306–319 (2000).
- 3. Da Costa, N. M. & Martin, K. A. C. The proportion of synapses formed by the axons of the lateral geniculate nucleus in layer 4 of area 17 of the cat. *J. Comp. Neurol.* **516**, 264–276 (2009).
- Markov, N. T. et al. Weight consistency specifies regularities of macaque cortical networks. Cereb. Cortex 21, 1254–1272 (2011).
- Bruno, R. M. & Sakmann, B. Cortex is driven by weak but synchronously active thalamocortical synapses. Science 312, 1622–1627 (2006).
- 6. Swindale, N. V. Visual map. Scholarpedia J. 3, 4607 (2008).
- Hubel, D. H. & Wiesel, T. N. Receptive fields of single neurones in the cat's striate cortex. J. Physiol. 148, 574–591 (1959).
- Riesenhuber, M. & Poggio, T. Hierarchical models of object recognition in cortex. Nat. Neurosci. 2, 1019–1025 (1999).
- 9. Angelucci, A. et al. Circuits for local and global signal integration in primary visual cortex. *J. Neurosci.* **22**, 8633–8646 (2002).
- Douglas, R. J., Koch, C., Mahowald, M., Martin, K. A. & Suarez, H. H. Recurrent excitation in neocortical circuits. *Science* 269, 981–985 (1995).
- Sompolinsky, H. & Shapley, R. New perspectives on the mechanisms for orientation selectivity. *Curr. Opin. Neurobiol.* 7, 514–522 (1997).
- Blakemore, C. & Tobin, E. A. Lateral inhibition between orientation detectors in the cat's visual cortex. Exp. Brain Res. 15, 439–440 (1972).
- Allman, J., Miezin, F. & McGuinness, E. Stimulus specific responses from beyond the classical receptive field: neurophysiological mechanisms for local-global comparisons in visual neurons. *Annu. Rev. Neurosci.* 8, 407–430 (1985).
- 14. Field, D. J., Hayes, A. & Hess, R. F. Contour integration by the human visual system: evidence for a local 'association field'. *Vis. Res.* **33**, 173–193 (1993).
- Gilbert, C. D. Adult cortical dynamics. *Physiol. Rev.* 78, 467–485 (1998).
- Albright, T. D. & Stoner, G. R. Contextual influences on visual processing. Annu. Rev. Neurosci. 25, 339–379 (2002).
- Kietzmann, T. C. et al. Recurrence is required to capture the representational dynamics of the human visual system. *Proc. Natl Acad. Sci. USA* 116, 21854–21863 (2019).
- Kar, K., Kubilius, J., Schmidt, K., Issa, E. B. & DiCarlo, J. J. Evidence that recurrent circuits are critical to the ventral stream's execution of core object recognition behavior. *Nat. Neurosci.* 22, 974–983 (2019).
- Bringuier, V., Chavane, F., Glaeser, L. & Frégnac, Y. Horizontal propagation of visual activity in the synaptic integration field of area 17 neurons. Science 283, 695–699 (1999).
- Roland, P. E. et al. Cortical feedback depolarization waves: a mechanism of top-down influence on early visual areas. *Proc. Natl Acad. Sci. USA* 103, 12586–12591 (2006).
- Xu, W., Huang, X., Takagaki, K. & Wu, J.-Y. Compression and reflection of visually evoked cortical waves. *Neuron* 55, 119–129 (2007).
- Nauhaus, I., Busse, L., Carandini, M. & Ringach, D. L. Stimulus contrast modulates functional connectivity in visual cortex. *Nat. Neurosci.* 12, 70–76 (2009).
- Reimer, A., Hubka, P., Engel, A. K. & Kral, A. Fast propagating waves within the rodent auditory cortex. Cereb. Cortex 21, 166–177 (2011).

- 24. Townsend, R. G. et al. Emergence of complex wave patterns in primate cerebral cortex. *J. Neurosci.* **35**, 4657–4662 (2015).
- Slovin, H., Arieli, A., Hildesheim, R. & Grinvald, A. Long-term voltage-sensitive dye imaging reveals cortical dynamics in behaving monkeys. J. Neurophysiol. 88, 3421–3438 (2002).
- Sato, T. K., Nauhaus, I. & Carandini, M. Traveling waves in visual cortex. Neuron 75, 218–229 (2012).
- 27. Davis, Z. W., Muller, L., Martinez-Trujillo, J., Sejnowski, T. & Reynolds, J. H. Spontaneous travelling cortical waves gate perception in behaving primates. *Nature* **587**, 432–436 (2020).
- Muller, L., Reynaud, A., Chavane, F. & Destexhe, A. The stimulusevoked population response in visual cortex of awake monkey is a propagating wave. *Nat. Commun.* 5, 3675 (2014).
- Davis, Z. et al. Spontaneous traveling waves naturally emerge from horizontal fiber time delays and travel through locally asynchronous-irregular states. *Nat. Commun.* 12, 6057 (2021).
- 30. Takahashi, K. et al. Large-scale spatiotemporal spike patterning consistent with wave propagation in motor cortex. *Nat. Commun.* **6**, 7169 (2015).
- 31. Schaffer, E. S., Ostojic, S. & Abbott, L. F. A complex-valued firingrate model that approximates the dynamics of spiking networks. *PLoS Comput. Biol.* **9**, e1003301 (2013).
- 32. Bassey, J., Qian, L. & Li, X. A survey of complex-valued neural networks. Preprint at https://arxiv.org/abs/2101.12249 (2021).
- 33. Trabelsi, C. et al. Deep complex networks. *International Conference on Learning Representations* (2018).
- Heeger, D. J. & Mackey, W. E. Oscillatory recurrent gated neural integrator circuits (ORGaNICs), a unifying theoretical framework for neural dynamics. *Proc. Natl Acad. Sci. USA* 116, 22783–22794 (2019).
- Krizhevsky, A., Sutskever, I. & Hinton, G. E. ImageNet classification with deep convolutional neural networks. *Commun. ACM* 60, 84–90 (2017).
- Budzinski, R. C. et al. Geometry unites synchrony, chimeras, and waves in nonlinear oscillator networks. Chaos 32, 031104 (2022).
- 37. Muller, L., Chavane, F., Reynolds, J. & Sejnowski, T. J. Cortical travelling waves: mechanisms and computational principles. *Nat. Rev. Neurosci.* **19**, 255–268 (2018).
- 38. Hellwig, B. A quantitative analysis of the local connectivity between pyramidal neurons in layers 2/3 of the rat visual cortex. *Biol. Cybern.* **82**, 111–121 (2000).
- Binzegger, T., Douglas, R. J. & Martin, K. A. C. A quantitative map of the circuit of cat primary visual cortex. *J. Neurosci.* 24, 8441–8453 (2004).
- Girard, P., Hupé, J. M. & Bullier, J. Feedforward and feedback connections between areas V1 and V2 of the monkey have similar rapid conduction velocities. J. Neurophysiol. 85, 1328–1331 (2001).
- 41. Elman, J. L. Finding structure in time. Cogn. Sci. 14, 179-211 (1990).
- 42. Jaeger, H. & Haas, H. Harnessing nonlinearity: predicting chaotic systems and saving energy in wireless communication. *Science* **304**, 78–80 (2004).
- 43. Pathak, J., Lu, Z., Hunt, B. R., Girvan, M. & Ott, E. Using machine learning to replicate chaotic attractors and calculate Lyapunov exponents from data. *Chaos* **27**, 121102 (2017).
- Vlachas, P. R. et al. Backpropagation algorithms and Reservoir Computing in Recurrent Neural Networks for the forecasting of complex spatiotemporal dynamics. *Neural Netw.* 126, 191–217 (2020).
- 45. Pathak, J., Hunt, B., Girvan, M., Lu, Z. & Ott, E. Model-free prediction of large spatiotemporally chaotic systems from data: a reservoir computing approach. *Phys. Rev. Lett.* **120**, 024102 (2018).
- Tang, Y., Kurths, J., Lin, W., Ott, E. & Kocarev, L. Introduction to focus issue: when machine learning meets complex systems: networks, chaos, and nonlinear dynamics. *Chaos* 30, 063151 (2020).

- Mathieu, M., Couprie, C. & LeCun, Y. Deep multi-scale video prediction beyond mean square error. *International Conference on Learning Representations* (2016).
- 48. Villegas, R., Yang, J., Hong, S., Lin, X. & Lee, H. Decomposing motion and content for natural video sequence prediction. *International Conference on Learning Representations* (2017).
- Desai, P. et al. Next frame prediction using ConvLSTM. J. Phys. Conf. Ser. 2161, 012024 (2022).
- 50. Michalski, V., Memisevic, R. & Konda, K. Modeling deep temporal dependencies with recurrent grammar cells. *Advances in Neural Information Processing Systems* **27**, (2014)
- Lotter, W., Kreiman, G. & Cox, D. Deep predictive coding networks for video prediction and unsupervised learning. *International Con*ference on Learning Representations (2017).
- Choi, M. & Tani, J. Predictive coding for dynamic visual processing: development of functional hierarchy in a multiple spatiotemporal scales RNN model. *Neural Comput.* 30, 237–270 (2018).
- Kwon, Y.-H. & Park, M.-G. Predicting future frames using retrospective cycle GAN. In Proc. 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 1811–1820 (IEEE, 2019).
- 54. Shouno, O. Photo-realistic video prediction on natural videos of largely changing frames. Preprint at https://arxiv.org/abs/2003.08635 (2020).
- Yu, W., Lu, Y., Easterbrook, S. & Fidler, S. Efficient and informationpreserving future frame prediction and beyond. In *Proc. 2020 International Conference on Learning Representations*. https:// openreview.net/pdf?id=F4e26c-K1DM (ICLR, 2020).
- Kasaraneni, S. H. Autoencoding video latents for adversarial video generation. Preprint at https://arxiv.org/abs/2201.06888 (2022).
- Ranzato, M. et al. Video (language) modeling: a baseline for generative models of natural videos. Preprint at https://arxiv.org/abs/1412.6604 (2014).
- Hou, R., Chang, H., Ma, B. & Chen, X. Video prediction with bidirectional constraint network. In Proc. 2019 14th IEEE International Conference on Automatic Face & Gesture Recognition (FG 2019), 1–8 (2019).
- 59. Wang, Z., Bovik, A. C., Sheikh, H. R. & Simoncelli, E. P. Image quality assessment: from error visibility to structural similarity. *IEEE Trans. Image Process.* **13**, 600–612 (2004).
- Ferster, D., Chung, S. & Wheat, H. Orientation selectivity of thalamic input to simple cells of cat visual cortex. *Nature* 380, 249–252 (1996).
- Ferster, D. & Miller, K. D. Neural mechanisms of orientation selectivity in the visual cortex. Annu. Rev. Neurosci. 23, 441–471 (2000).
- Stettler, D. D., Das, A., Bennett, J. & Gilbert, C. D. Lateral connectivity and contextual interactions in macaque primary visual cortex. *Neuron* 36, 739–750 (2002).
- Moshe, B., Lena, G., Eli, S., Michal, I. & Ronen, B. Actions as spacetime shapes. in Proc. Tenth IEEE International Conference on Computer Vision, Beijing, China, 17–20 (2005).
- 64. Thorpe, S., Fize, D. & Marlot, C. Speed of processing in the human visual system. *Nature* **381**, 520–522 (1996).
- 65. McBeath, M. K. The rising fastball: baseball's impossible pitch. *Perception* **19**, 545–552 (1990).
- Chelazzi, L., Miller, E. K., Duncan, J. & Desimone, R. A neural basis for visual search in inferior temporal cortex. *Nature* 363, 345–347 (1993).
- Wheeler, M. E., Petersen, S. E. & Buckner, R. L. Memory's echo: vivid remembering reactivates sensory-specific cortex. *Proc. Natl Acad.* Sci. USA 97, 11125–11129 (2000).
- 68. Horner, A. J., Bisby, J. A., Bush, D., Lin, W.-J. & Burgess, N. Evidence for holistic episodic recollection via hippocampal pattern completion. *Nat. Commun.* **6**, 7462 (2015).
- Vanni, S., Hokkanen, H., Werner, F. & Angelucci, A. Anatomy and physiology of macaque visual cortical areas V1, V2, and V5/MT:

- bases for biologically realistic models. *Cereb. Cortex* **30**, 3483–3517 (2020).
- Benvenuti, G. et al. Anticipatory responses along motion trajectories in awake monkey area V1. Preprint at bioRxiv https://doi.org/10.1101/2020.03.26.010017 (2020).
- 71. Montbrió, E., Pazó, D. & Roxin, A. Macroscopic description for networks of spiking neurons. *Phys. Rev. X* **5**, 021028 (2015).
- Zanos, T. P., Mineault, P. J., Nasiotis, K. T., Guitton, D. & Pack, C. C. A sensorimotor role for traveling waves in primate visual cortex. Neuron 85, 615–627 (2015).
- Townsend, R., Solomon, S. S., Martin, P. R., Solomon, S. G. & Gong, P. Visual motion discrimination by propagating patterns in primate cerebral cortex. *J. Neurosci.* 37, 10074–10084 (2017).
- Aggarwal, A. et al. Visual evoked feedforward-feedback traveling waves organize neural activity across the cortical hierarchy in mice. Nat. Commun. 13, 1-16 (2022).
- LeCun, Y., Bengio, Y. & Hinton, G. Deep learning. *Nature* 521, 436–444 (2015).
- Xu, Y., Xiao, T., Zhang, J., Yang, K. & Zhang, Z. Scale-invariant convolutional neural networks. Preprint at https://arxiv.org/abs/ 1411.6369 (2014).
- 77. Hess, R. F. & Dakin, S. C. Contour integration in the peripheral field. *Vis. Res.* **39**, 947–959 (1999).
- Olshausen, B. A. How close are we to understanding V1? Neural Comput. 17, 1665–1699 (2005).
- Schrimpf, M. et al. Brain-score: which artificial neural network for object recognition is most brain-like? Preprint at bioRxiv https://doi. org/10.1101/407007 (2020).
- Cadieu, C. F. & Olshausen, B. A. Learning intermediate-level representations of form and motion from natural movies. *Neural Comput.* 24, 827–866 (2012).
- Acebrón, J. A., Bonilla, L. L., Pérez Vicente, C. J., Ritort, F. & Spigler, R. The Kuramoto model: a simple paradigm for synchronization phenomena. Rev. Mod. Phys. 77, 137–185 (2005).
- 82. Rodrigues, F. A., Peron, T. K. D. M., Ji, P. & Kurths, J. The Kuramoto model in complex networks. *Phys. Rep.* **610**, 1–98 (2016).
- 83. Breakspear, M., Heitmann, S. & Daffertshofer, A. Generative models of cortical oscillations: neurobiological implications of the Kuramoto model. *Front. Hum. Neurosci.* **4**, 190 (2010).
- 84. Muller, L., Mináč, J. & Nguyen, T. T. Algebraic approach to the Kuramoto model. *Phys. Rev. E* **104**, L022201 (2021).
- Heim, N. & Avery, J. E. Adaptive anomaly detection in chaotic time series with a spatially aware echo state network. Preprint at https:// arxiv.org/abs/1909.01709 (2019).
- Gorelick, L., Blank, M., Shechtman, E., Irani, M. & Basri, R. Actions as space-time shapes. *IEEE Trans. Pattern Anal. Mach. Intell.* 29, 2247–2253 (2007).
- 87. Anselmi, F. et al. Unsupervised learning of invariant representations. *Theor. Comput. Sci.* **633**, 112–121 (2016).
- Benigno, G. B. et al. Waves traveling over a map of visual space can ignite short-term predictions of sensory input. mullerlab/ benignoEAwavecomp: v1.0.0. Zenodo https://doi.org/10.5281/ zenodo.7863700 (2023).

#### Acknowledgements

This work was supported by the Canadian Institute for Health Research (L.M.) and NSF (NeuroNex Grant No. 2015276) (L.M.), NIH Grant R01-EY028723 (J.H.R.), BrainsCAN at Western University through the Canada First Research Excellence Fund (CFREF) (L.M.), the Fiona and Sanjay Jha Chair in Neuroscience (J.H.R.), Compute Ontario (computeontario.ca) (L.M.), Digital Research Alliance of Canada (https://alliancecan.ca/en) (L.M.), SPIRITS 2020 of Kyoto University (L.M.), and the Western Academy for Advanced Research (L.M.). R.C.B. gratefully acknowledges the Western Institute for Neuroscience Clinical

Research Postdoctoral Fellowship. G.B.B. gratefully acknowledges the Canadian Open Neuroscience Platform (Graduate Scholarship), the Vector Institute (Postgraduate Affiliate), and the National Sciences and Engineering Research Council of Canada (Canada Graduate Scholarship—Doctoral). The authors thank Alex Busch for her help with illustrations.

#### **Author contributions**

Conceptualization: G.B.B., L.M.; data curation: G.B.B.; formal analysis: G.B.B., L.M.; funding acquisition: J.H.R., L.M.; investigation: G.B.B., L.M.; methodology: G.B.B., R.C.B., Z.W.D., J.H.R., L.M.; supervision: J.H.R., L.M.; visualization: G.B.B.; writing—original draft: G.B.B., L.M.; writing—review and editing: G.B.B., R.C.B., Z.W.D., J.H.R., L.M.

#### **Competing interests**

The authors declare no competing interests.

#### **Additional information**

**Supplementary information** The online version contains supplementary material available at https://doi.org/10.1038/s41467-023-39076-2.

**Correspondence** and requests for materials should be addressed to Lyle Muller.

**Peer review information** *Nature Communications* thanks Alex Proekt and the other anonymous reviewers for their contribution to the peer review of this work.

**Reprints and permissions information** is available at http://www.nature.com/reprints

**Publisher's note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this license, visit http://creativecommons.org/licenses/by/4.0/.

© The Author(s) 2023