

## MIT Open Access Articles

*Robust and brain-like working memory  
through short-term synaptic plasticity*

The MIT Faculty has made this article openly available. **Please share** how this access benefits you. Your story matters.

**Citation:** Kozachkov, Leo, Tauber, John, Lundqvist, Mikael, Brincat, Scott L, Slotine, Jean-Jacques et al. 2022. "Robust and brain-like working memory through short-term synaptic plasticity." PLOS Computational Biology, 18 (12).

**As Published:** 10.1371/journal.pcbi.1010776

**Publisher:** Public Library of Science (PLoS)

**Persistent URL:** <https://hdl.handle.net/1721.1/147183>

**Version:** Final published version: final published article, as it appeared in a journal, conference proceedings, or other formally published context

**Terms of use:** Creative Commons Attribution 4.0 International license



## RESEARCH ARTICLE

## Robust and brain-like working memory through short-term synaptic plasticity

Leo Kozachkov<sup>1,2,3‡</sup>, John Tauber<sup>1,2,4‡</sup>, Mikael Lundqvist<sup>1,2,5</sup>, Scott L. Brincat<sup>1</sup>, Jean-Jacques Slotine<sup>2,3</sup>, Earl K. Miller<sup>1,2\*</sup>

**1** The Picower Institute for Learning & Memory, Massachusetts Institute of Technology (MIT), Cambridge, Massachusetts, United States of America, **2** Department of Brain & Cognitive Sciences, Massachusetts Institute of Technology (MIT), Cambridge, Massachusetts, United States of America, **3** Nonlinear Systems Laboratory, Massachusetts Institute of Technology (MIT), Cambridge, Massachusetts, United States of America, **4** Institute for Data, Systems, and Society, Massachusetts Institute of Technology (MIT), Cambridge, Massachusetts, United States of America, **5** Department of Clinical Neuroscience, Karolinska Institute, Stockholm, Sweden

‡ These authors are co-first authors on this work.

\* [ekmiller@mit.edu](mailto:ekmiller@mit.edu)

## OPEN ACCESS

**Citation:** Kozachkov L, Tauber J, Lundqvist M, Brincat SL, Slotine J-J, Miller EK (2022) Robust and brain-like working memory through short-term synaptic plasticity. *PLoS Comput Biol* 18(12): e1010776. <https://doi.org/10.1371/journal.pcbi.1010776>

**Editor:** Blake A. Richards, McGill University, CANADA

**Received:** February 2, 2022

**Accepted:** November 29, 2022

**Published:** December 27, 2022

**Copyright:** © 2022 Kozachkov et al. This is an open access article distributed under the terms of the [Creative Commons Attribution License](https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

**Data Availability Statement:** The code and data are available here: [https://github.com/kozleol/robust\\_wm\\_stsp](https://github.com/kozleol/robust_wm_stsp).

**Funding:** This work was supported by Office of Naval Research N00014-22-1-2453 (E.K.M.), The JPB Foundation (E.K.M.), ERC Starting Grant 949131 (M.L.), and VR Starting Grant 2018-04197 (M.L.). The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

## Abstract

Working memory has long been thought to arise from sustained spiking/ attractor dynamics. However, recent work has suggested that short-term synaptic plasticity (STSP) may help maintain attractor states over gaps in time with little or no spiking. To determine if STSP endows additional functional advantages, we trained artificial recurrent neural networks (RNNs) with and without STSP to perform an object working memory task. We found that RNNs with and without STSP were able to maintain memories despite distractors presented in the middle of the memory delay. However, RNNs with STSP showed activity that was similar to that seen in the cortex of a non-human primate (NHP) performing the same task. By contrast, RNNs without STSP showed activity that was less brain-like. Further, RNNs with STSP were more robust to network degradation than RNNs without STSP. These results show that STSP can not only help maintain working memories, it also makes neural networks more robust and brain-like.

## Author summary

Working memory has been thought to depend on sustained spiking alone. But recent evidence shows that spiking is often sparse, not sustained. Short-term synaptic plasticity (STSP) could help by maintaining memories between spiking. To test this, we compared artificial recurrent neural networks (RNNs) with and without short-term synaptic plasticity (STSP). Both types of RNNs could maintain working memories. But RNNs with STSP functioned better. They were more robust to network degradation. Plus, their activity was more brain-like than RNNs without STSP. These results support a role for STSP in working memory.

**Competing interests:** The authors have no competing interests.

## Introduction

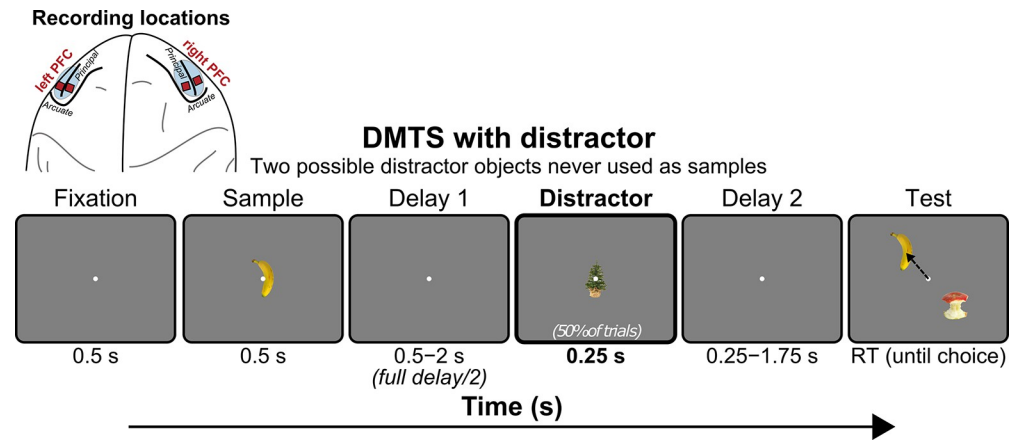
Working memory (WM), the holding of information “online” and available for processing, is central to higher cognitive functions [1,2]. A well-established neural correlate of WM is spiking over a memory delay [3–5]. For many years, this was thought to be the sole mechanism underlying WM maintenance. The idea is that sensory inputs elicit unique patterns of spiking that are sustained via recurrent connections [6], creating attractor states—stable patterns of activity that retain the WM [7]. It seems evident that these attractor dynamics play an important role in WM. Recent observations, however, have suggested that there may be more going on [8–11]. A few neurons seem to show spiking that looks persistent enough to be an attractor state, but the bulk of neurons show memory delay spiking that is sparse [12–15]. This is especially true when spiking is examined in real time (i.e., on single trials) because averaging across trials can create the appearance of persistence even when the underlying activity is quite sparse.

This all begs the question of how WMs are maintained over these gaps in time with little-to-no spiking [11,16]. One possibility was suggested by observations of short-term synaptic plasticity (STSP), transient (< 1 second) changes in synaptic weights induced by spiking, in circuits in the prefrontal cortex [17]. Several groups have suggested updating the attractor dynamics model with this feature [18–20]. The idea is that STSP helps the spiking. Spikes induce a transient “impression” in the synaptic weights that can maintain the network state between spikes [21–23]. Evidence for STSP comes from techniques like patch-clamp recording that are difficult to implement in the working brain, especially in NHPs. Thus, we tested the role of STSP in WM by using computational modeling in conjunction with “ground-truthing” via analysis of spiking recorded from the PFC of a NHP performing a WM task. We found that PFC spiking carried little-to-no stimulus-specific WM information across the delay. We aimed to determine if network models with STSP can solve the working memory task, whether they have properties similar to those seen in the actual brain, and whether STSP endows functional advantages. The answer to these questions was “yes”.

We trained Recurrent Neural Networks (RNNs) with and without STSP to test how it affects network performance and function. We focused on the key property of robustness [9,24–27]. Working memories must be maintained in the face of distractions. Networks need to deal with noise and show graceful degradation (i.e., continue to function when portions of the network are damaged). Our analysis showed RNNs with and without STSP were robust against distractors. However, *only the RNNs with STSP were “brain-like”*—their activity more closely resembled activity recorded from the prefrontal cortex of a NHP performing a WM task. RNNs with STSP were also more robust against synaptic ablation. Thus, STSP offers functional advantages and explains how WM can be maintained between stimulus presentations. This point has been made repeatedly in the literature. However, our study is the first to train artificial neural networks with STSP and quantitatively measure their similarity to recorded electrophysiological data.

## Results

A NHP was trained to perform an object delayed-match-to-sample task (Fig 1). The NHP was shown a sample object and had to choose its match after a variable-length memory delay. At mid-delay a distractor object (1 of 2 possible objects never used as samples) was presented (for 0.25s) on 50% randomly chosen trials. We recorded multi-unit activity (MUA) bilaterally in dorsolateral PFC (dlPFC) and ventrolateral PFC (vlPFC) using four 64-electrode Utah arrays for a total of 256 electrodes. The animal learned to do the task consistently at ~99% accuracy for both distractor and non-distractor trials.



**Fig 1. Electrode location and task structure.** Utah arrays were implanted bilaterally in dorsolateral PFC (dlPFC) and ventrolateral PFC (vlPFC). Animal performed a distracted delayed match-to-sample task. Each trial began with visual fixation on the middle of the screen for 0.5s. Fixation was maintained throughout the trial until the behavioral response. The delay length was parametrically varied from 1–4 s in five logarithmic steps, randomly chosen each trial. At mid-delay a neutral distractor (1 of 2 possible objects never used as samples) was presented randomly on 50% of trials. During the multi-choice test the NHP was allowed to freely saccade between all objects on the screen. The final choice was indicated by fixating on it for at least one second.

<https://doi.org/10.1371/journal.pcbi.1010776.g001>

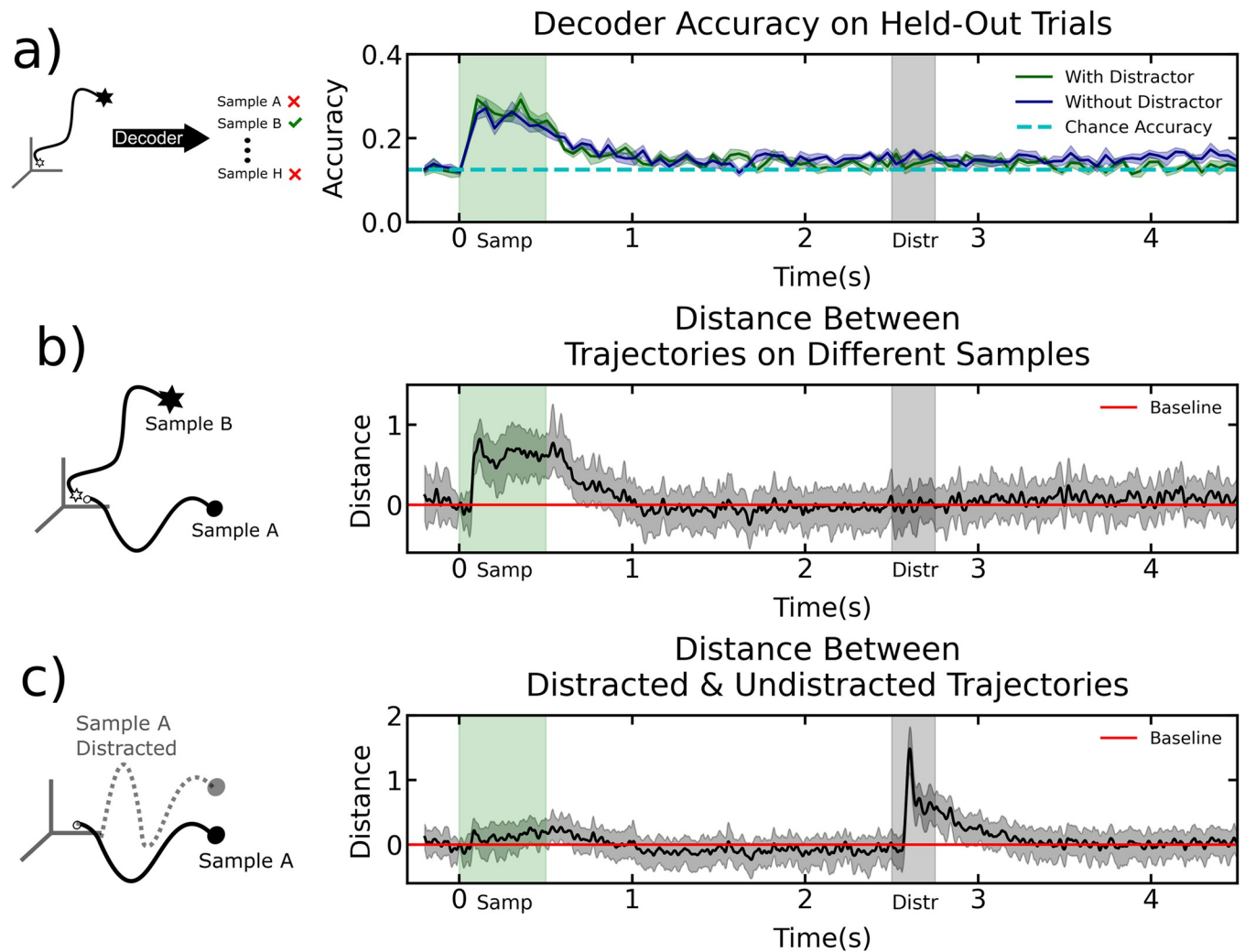
### Sample information in population neural activity was weak over longer delays

First, we examined MUA recorded from the lateral PFC. To quantify the amount of sample object information carried by spiking, we used a linear classifier (see [Methods](#) for details)<sup>1</sup>. This showed that spiking carried sample object information for about one second after the sample disappeared. From the start of the delay period the decoder accuracy decreased steadily towards chance (Figs 2A and S1). We corroborated this by measuring the distance between neural population activity for all pairs of sample objects. That gave an average distance between experimental conditions at every timepoint (see [Methods](#)). This showed that the distance between population MUA activity for different samples returned to pre-stimulus levels (Figs 2B and S2). Interestingly, we found that this was *not* simply due to spiking returning to pre-sample values. We determined this by training a classifier to discriminate between pre and post sample spiking activity. We found that this classifier was consistently able to discriminate between pre and post sample spiking activity over the delay (S3 Fig).

### Population neural activity was robust to distractors

We found that when the mid-delay distractor was presented, neural trajectories diverged from that of non-distractor trials (Figs 2C, S4). Once the distractor disappeared, the trajectories quickly reconvened, indicating neural stability. This was true for all five delay lengths used (1–4 s in five logarithmic steps, S4 Fig). The time course of trajectory reconvening was roughly exponential with a time constant of ~200 milliseconds (Figs 2C, S4). We determined this by fitting an exponential function to the recovery curves and measuring the inverse of the fitted decay constants. This is consistent with prior observations that time constants in cortex peak at this value in the PFC [28].

This all raises a couple of questions. 1. How can PFC networks support WM task performance when, across the population, sample information in spiking is relatively weak? 2. How do PFC networks achieve the stability to recover from distractors? To answer these questions,



**Fig 2.** a) Left: training a decoder to predict sample identity given a neural trajectory. Right: decoder accuracy on held-out trials for distracted vs. undistracted trials b) Left: comparing trial-averaged trajectories corresponding to different samples. Right: average pairwise distance in state-space between trajectories elicited by all possible sample images. Normalized by the average pre-stim distance. c) Left: comparing trial-averaged distracted vs. non-distracted trajectories through neural state space. Right: distance in state space between distracted vs. non-distracted trajectories throughout the trial. Shown are trials with a delay of four seconds.

<https://doi.org/10.1371/journal.pcbi.1010776.g002>

we used RNN modeling and neural network theory. As we will show, the two questions share a common answer: STSP.

### RNNs with STSP are more brain-like

RNNs with and without STSP were able to successfully perform the object delayed match to sample task. However, only the RNNs with STSP did so with activity that was similar to that seen in the actual PFC. Our main hypothesis space consisted of four different kinds of RNNs: two with fixed synaptic weights (fixed after training) and two with STSP. The two fixed-weight networks were ‘vanilla’ RNNs. They differed only in their activation functions. One used a hyperbolic tan ( $\tanh$ ), the other used a rectified linear (ReLU). We will refer to these fixed synapse networks as *FS-tanh* and *FS-relu* respectively. We choose these two activation functions because they represent two sides of a spectrum. Tanh units can become unresponsive for very

large inputs (i.e. saturate) while ReLU units cannot. Both activations are commonly used throughout computational neuroscience and machine learning [29,30]. They lead RNNs to prefer one strategy over another for performing a task. ReLU units are ideal for forming line attractors, while tanh units are ideal for forming point attractors [30]. Additionally, we also trained Long-Short-Term-Memory networks [31] (LSTMs) and Gated Recurrent Unit networks [32] (GRUs). These two networks are popular in machine learning and are well-known for their excellent performance capabilities. However, they are not biologically-plausible models of the brain [33], and thus serve as a useful tool for dissociating our measures of brain-similarity and robustness.

The synaptic weights for the fixed-synapse networks are only adjusted during training. After training, they are untouched. The other two models use STSP to adjust synaptic weights during each trial. This represents the distinction between long-term and short-term memory. Long-term memory is acquired over the course of task-optimization and remains fixed throughout the trial. Short-term memory changes during the trial and reflects the contents of working memory. We reset the state of the plastic synapses at the beginning of each trial.

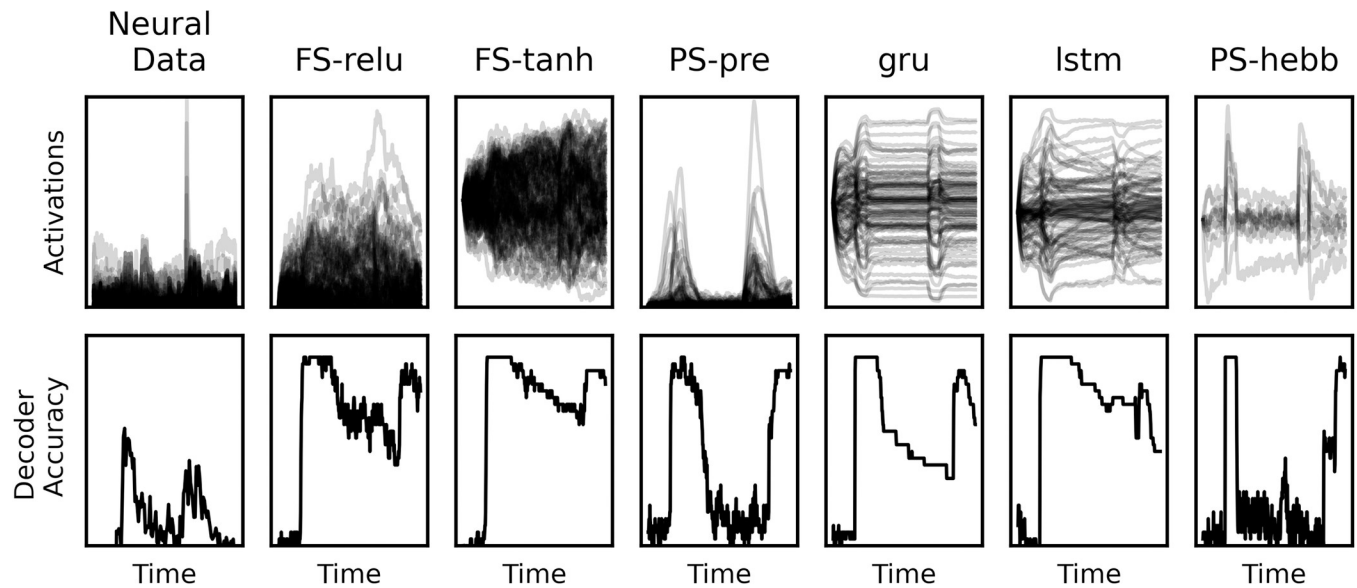
One RNN with STSP was based on a model introduced by Mongillo et al [19]. The model uses a set of simplified equations for synaptic calcium dynamics to endow the RNN with STSP. The Mongillo model adjusts synaptic weights based on presynaptic neural activity. For this reason we will call it *PS-pre*. The other RNN with STSP was introduced by Kozachkov et al. [34] and adjusts synaptic weights in a way that depends on both the pre and post synaptic neural activity. For this reason, we will call it *PS-hebb*. It uses excitatory anti-Hebbian and inhibitory Hebbian mechanisms to stabilize the RNN [34].

All models were trained with backpropagation-through-time using a standard deep learning library [35]. While the FS models were trained without any explicit constraints, the PS models needed to be parameterized to ensure they satisfy certain properties throughout training. In particular, the *PS-pre* model had separate excitatory and inhibitory neurons. The weight matrix therefore needed to be parameterized so that these populations always remained separate (see [Methods](#)). Likewise, for *PS-hebb* the weights were parameterized so the network remained stable throughout training. To ensure that our results did not depend on the particular choice of training hyperparameters, we trained the six models under a wide range of hyperparameter settings. These hyperparameters were: the degree of activity regularization, degree of parameter regularization during training, and the size of the network (called hidden size, referring to the number of hidden units). Roughly 2000 models were trained in total. We refer the reader to the [Methods](#) section for more details on the models and the training process. As in the analysis of the experimental data, we used a decoder to read out sample object information over time. We did this for both the fixed-synapse and plastic-synapse models. For the fixed-synapse RNNs we trained the decoder on trajectories from the spike rates. For the STSP RNNs, we trained two separate decoders: one on the spike-rates and one on the synaptic weights. In the LSTM and GRU models, there are no variables which can clearly be denoted as “STSP weights” (since they are not biological models), so we exclude them from this analysis.

In all models, sample presentation elicited a large increase in decoding accuracy (Figs 3 and 4). Spike rate decoding of sample information in the fixed-synapse models remained high throughout the memory delay. This is in sharp contrast to the actual PFC MUA data. The PFC MUA showed a large drop in sample object decoding accuracy once the sample disappeared and especially for delays longer than one second. This was mirrored in the STSP models (Figs 3 and 4). However, sample decoding accuracy on the synaptic weights for the STSP models remained high over the entire delay, including longer delays (Fig 4).

This raises the question: do STSP models better capture the dynamics of the real PFC? We addressed this question by comparing the similarity of each of the ~2000 trained models to the





**Fig 3. Example neural activations and their corresponding decoder curves.** Top row: neural activity corresponding to a single trial condition. The leftmost panel is the actual neural data. The other panels are artificial neural networks, whose details are described in the main text. Bottom row: decoder accuracy curves corresponding to the neural activations in the top row.

<https://doi.org/10.1371/journal.pcbi.1010776.g003>

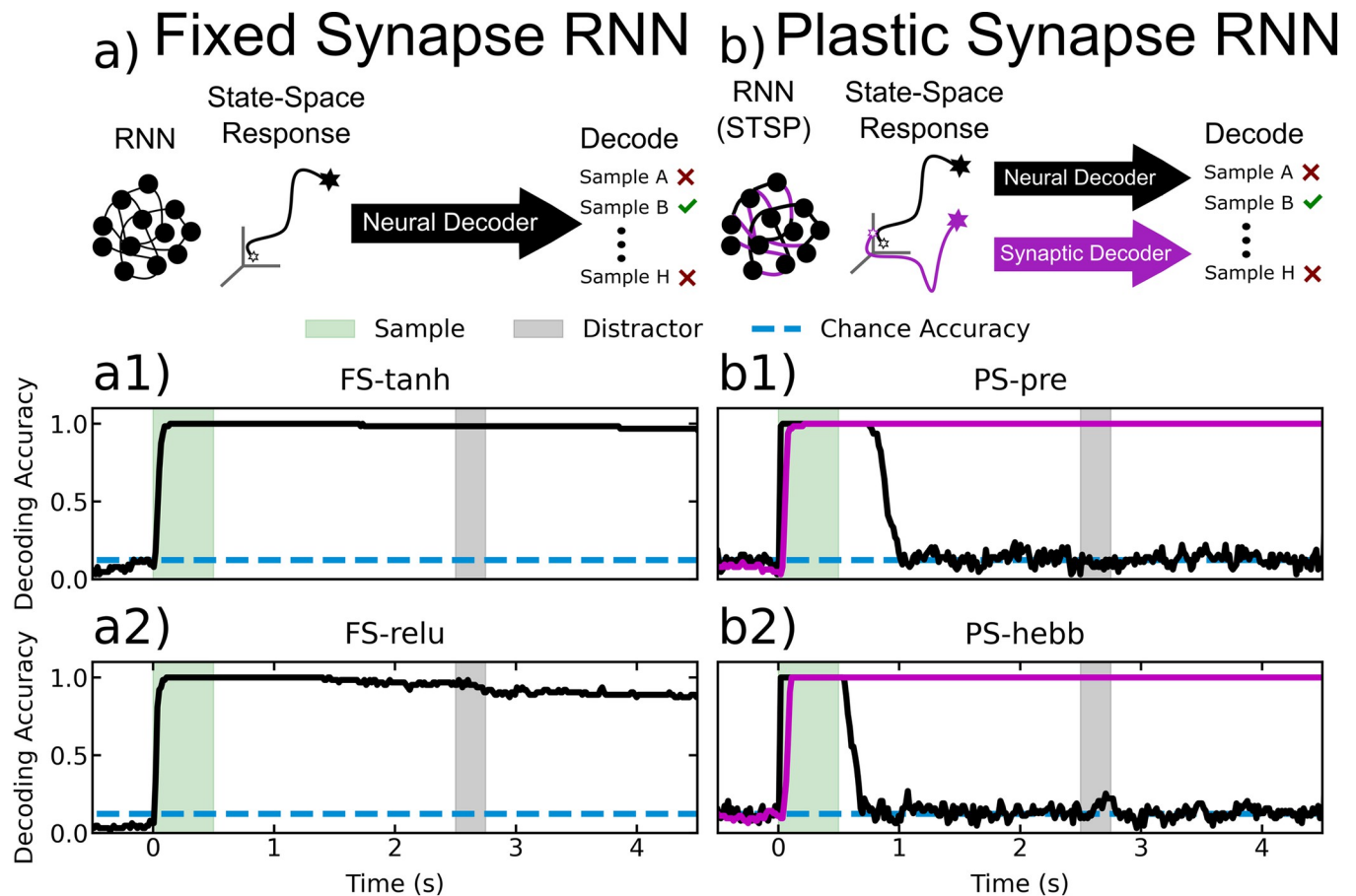
actual neural data. To quantify the similarity between model and brain, we used the Pearson correlation between the decoder accuracy curves for the real neural data and the RNNs (Fig 3, bottom row). We did this for each delay length and averaged the correlation scores to get a final brain-RNN correlation value for each model. We did this analysis with and without the LSTM and GRU networks, to dissociate the potential performance improvements of LSTM/GRU with their brain-similarity. The results without these networks are shown in Fig 5, the results with are shown in S8 Fig.

This analysis revealed that the *PS-hebb* model was the most similar to the brain across a wide range of hyperparameters (Fig 5, left column). Interestingly, for a particular choice of parameter regularization, this analysis also revealed that increasing the amount of activity regularization led to a corresponding increase in brain-similarity for the *PS-pre* model (Fig 5, top left panel). Within the fixed-synapse model class, the models that were most brain-like across all hyperparameter values showed a marked dip in neural decoder accuracy during the delay period—as expected. However, this dip was far less pronounced than in the models with plastic synapses. The result was that models without synaptic plasticity had a lower similarity to the actual neural data.

An important difference between measuring information in RNNs and brains is that we subsample neurons in the brain. It is possible that the decoding accuracy curves are affected by this subsampling. To control for this, we repeated the above analysis while subsampling at four different values. We randomly picked 0.01%, 1.0%, 10.0% and 100.0% of the neurons in each RNN to decode from. The result was that STSP networks were more brain-like and this did not change because of subsampling (S5 Fig).

### STSP increases structural robustness

The two models with STSP were almost always more *structurally* robust than the models with fixed synaptic weights (Fig 5, middle column). We examined how the performance accuracy of the trained network varied as we randomly ablated a varying fraction of synaptic weights. We



**Fig 4. Decoding accuracy for RNNs.** a) A decoder is trained to predict the sample label from RNN neural trajectories, for the fixed-synapse RNNs. a1) The accuracy of the trained decoder as a function of time from sample onset for FS-tanh. a2) The same plot as (a1), for FS-relu. b) Two decoders are trained on the neural and synaptic trajectories separately. Black lines indicate neural decoding, purple indicate synaptic decoding. b1) Neural and synaptic decoding accuracy as a function of time for PS-pre. b2) Same plot as in (b1) but for PS-hebb.

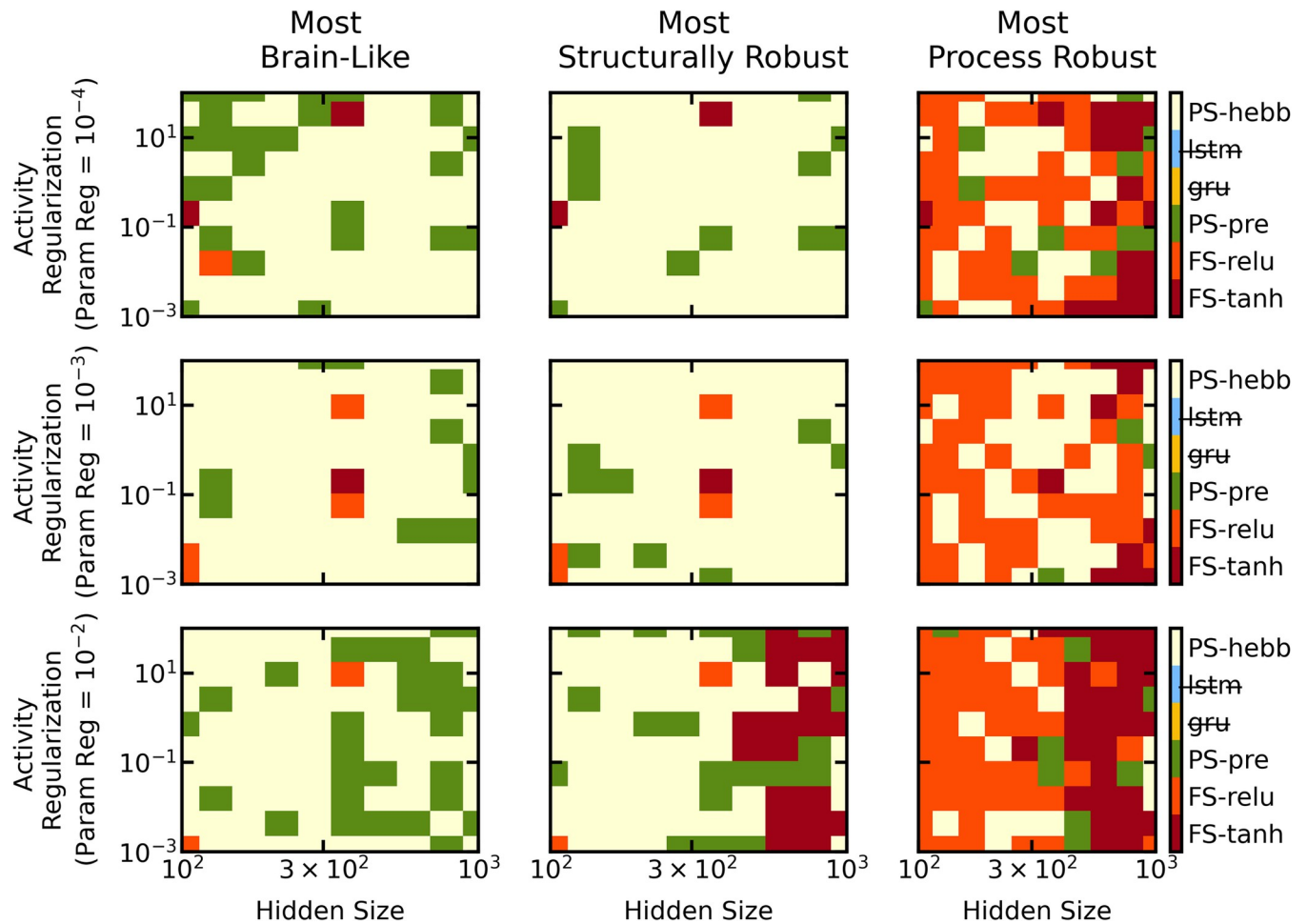
<https://doi.org/10.1371/journal.pcbi.1010776.g004>

found that the STSP network's performance on the task remained high even when as many as half of the synapses were ablated. By contrast, the fixed-synapse models were highly sensitive to synaptic ablation. Their performance severely degraded even with ablation of only 10–20% of the synapses.

We repeated this analysis over all the models trained with different hyperparameters. From each model we defined two measurements to quantify robustness. The first was robustness to structural noise. The second was the robustness to synaptic noise. Both measurements were calculated in a similar way. We took the weighted average RNN performance over all the noise values tested. We weighed each term in the sum by the corresponding noise value. The reason for this weighting scheme is simple. If the RNN performance is high when the noise is low, this does not tell us much about robustness. However, if the RNN performance is high when the noise is high, this does indicate robustness. Thus, a natural measure of robustness is the product of noise and performance.

This analysis revealed that the *PS-hebb* and *PS-pre* models were more robust to structural perturbations than the fixed synapse networks across a wide range of hyperparameters (Figs 5 and S8). Interestingly, the fixed synapse networks tended to be more robust to process noise (although always less brain-like). We repeated this analysis with the LSTM and GRU networks.



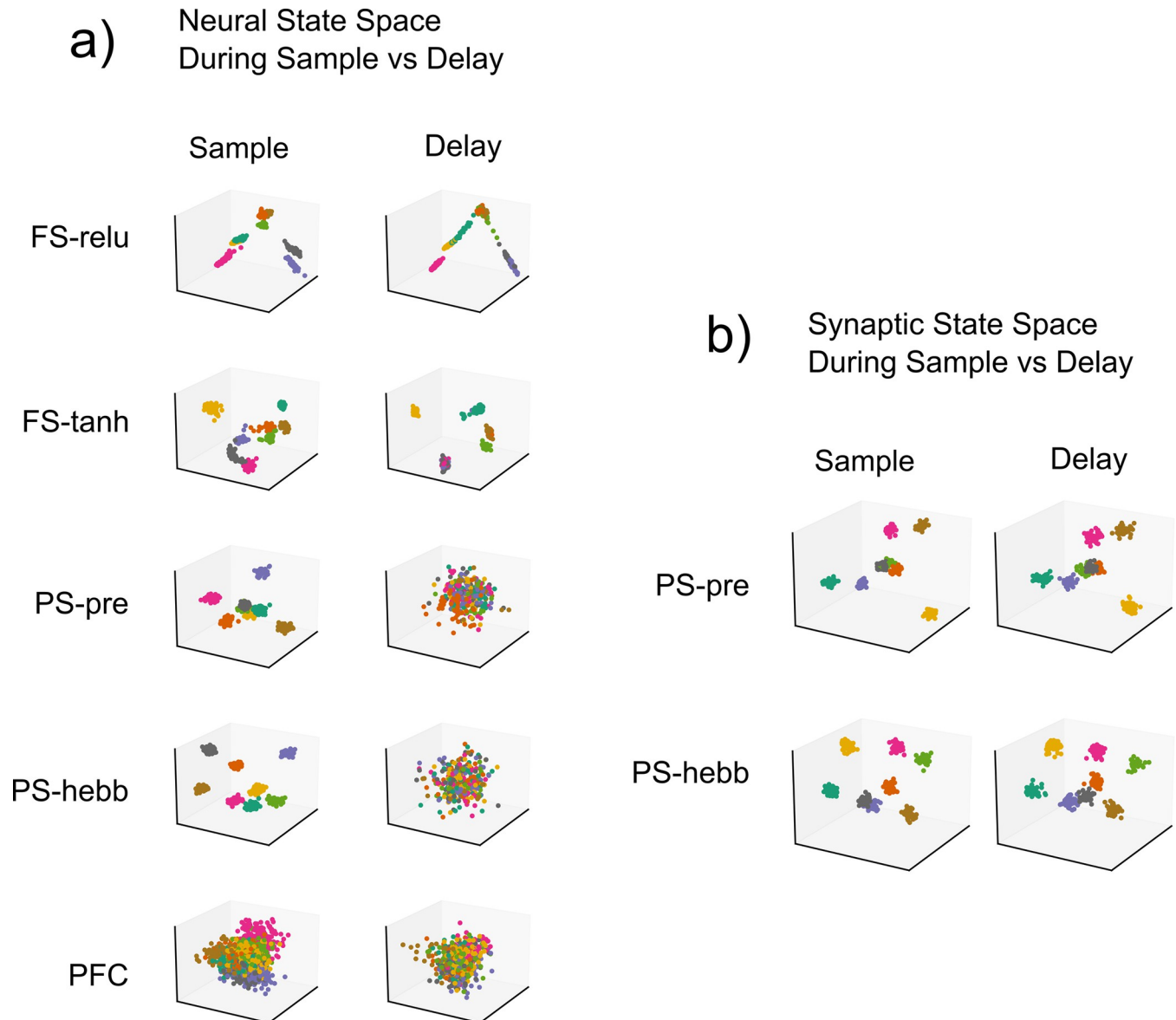


**Fig 5. Results of the hyperparameter sweep across number of hidden neurons, parameter regularization strength, and activity regularization strength.** Each row corresponds to a different parameter regularization ( $1e-4, 1e-3, 1e-2$ ). Each column corresponds to a different observed quantity (brain-likeness, structural robustness, process robustness). Each subplot is a  $10 \times 10$  grid, corresponding to 10 possible hidden size / activity regularization configurations. Shown in each square of that grid is the most brain-like/robust network corresponding to that particular hyperparameter configuration. LSTM and GRU networks were excluded (see S8 Fig for corresponding figure when they are included). Each color corresponds to a different network.

<https://doi.org/10.1371/journal.pcbi.1010776.g005>

This revealed that, out of the six kinds of networks, the LSTM and GRU networks were the most robust to structural noise and process noise across a wide range of hyperparameters (S8 Fig). However, they were almost always less brain-like than the *PS-hebb* and *PS-pre* models. This suggests that brain-similarity does not imply robustness, and vice versa. To better understand how the different trained models achieved WM that was robust to distractors we investigated how these networks organized their activity in state-space.

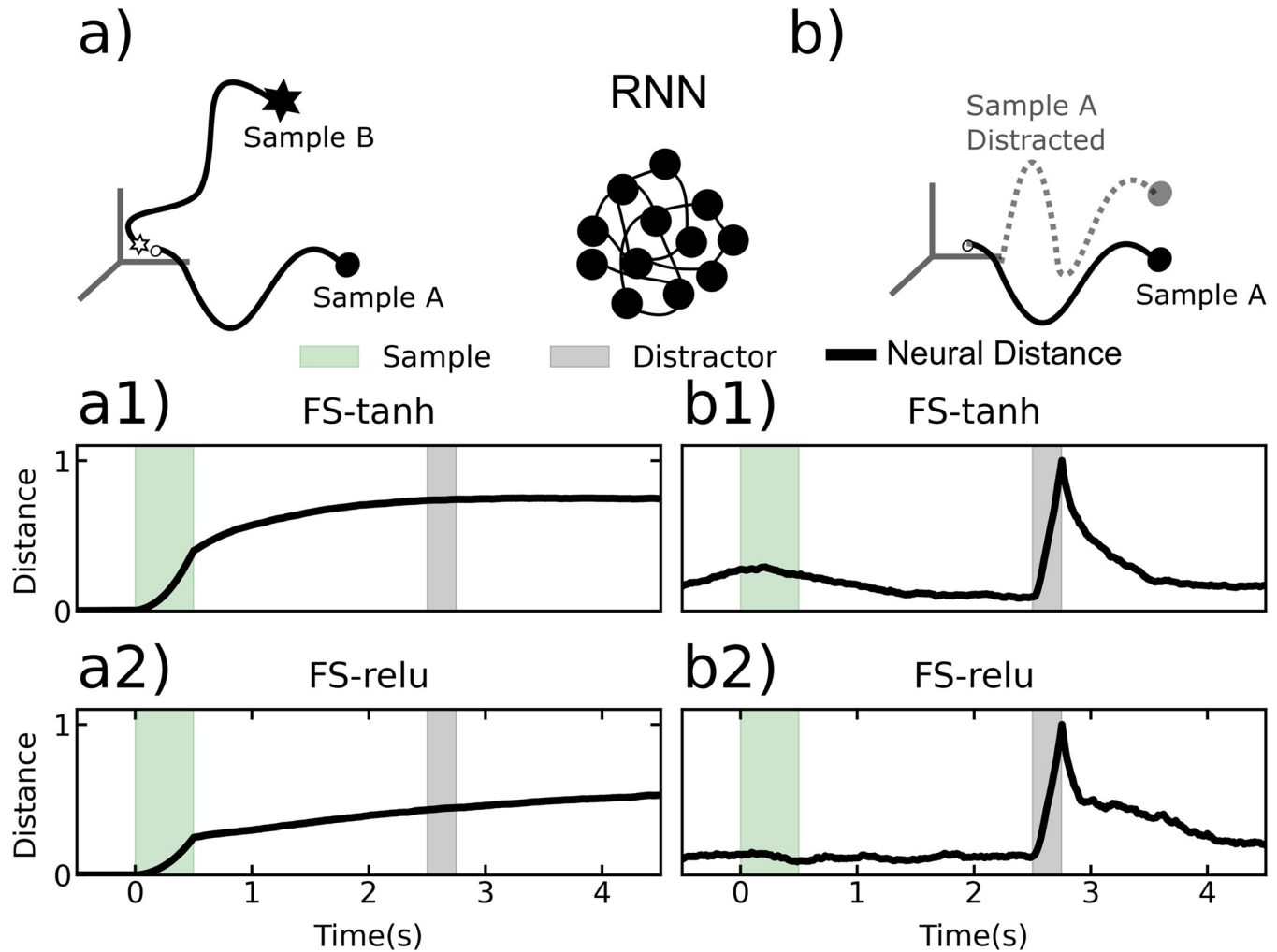
We found that the fixed-synapse RNNs learned to perform the task by using simple attractors. To visualize these attractors, we projected the high-dimensional RNN activity into a lower dimensional space using Linear Discriminant Analysis (Fig 6). We reasoned that this projection would give us the clearest visualization of the underlying state space attractors. For the plastic-synapse RNNs, we did this for both the neural and synaptic state-space. To better quantify the attractor properties of these networks, we measured the distances between state-space trajectories of the most brain-like models, as we did for the neural data (Fig 7). This revealed different trajectories for different sample objects. Each trajectory settled into a steady state that was unique for each sample, indicating the presence of attractors (Figs 6 and 7).



**Fig 6. Dimensionality reduced space plots for the most brain-like models.** Time-averaged RNN activity during the sample-period as well as 500ms before the end of the delay period. Linear Discriminant Analysis (LDA) was used to project the data into three-dimensions. To account for differences in training/test splits, an Orthogonal Procrustes operation was used to rotationally align the sample and delay period activity. Colors denote sample IDs. Fixed synapse models have activity organized around simple attractors in state space. There is one attractor for each sample ID. Plastic synapse models exhibit high sample-separability in synaptic state space, and limited separability in the neural state space during the delay period. Similarly, PFC exhibited higher neural discriminability during the sample than during the delay period.

<https://doi.org/10.1371/journal.pcbi.1010776.g006>

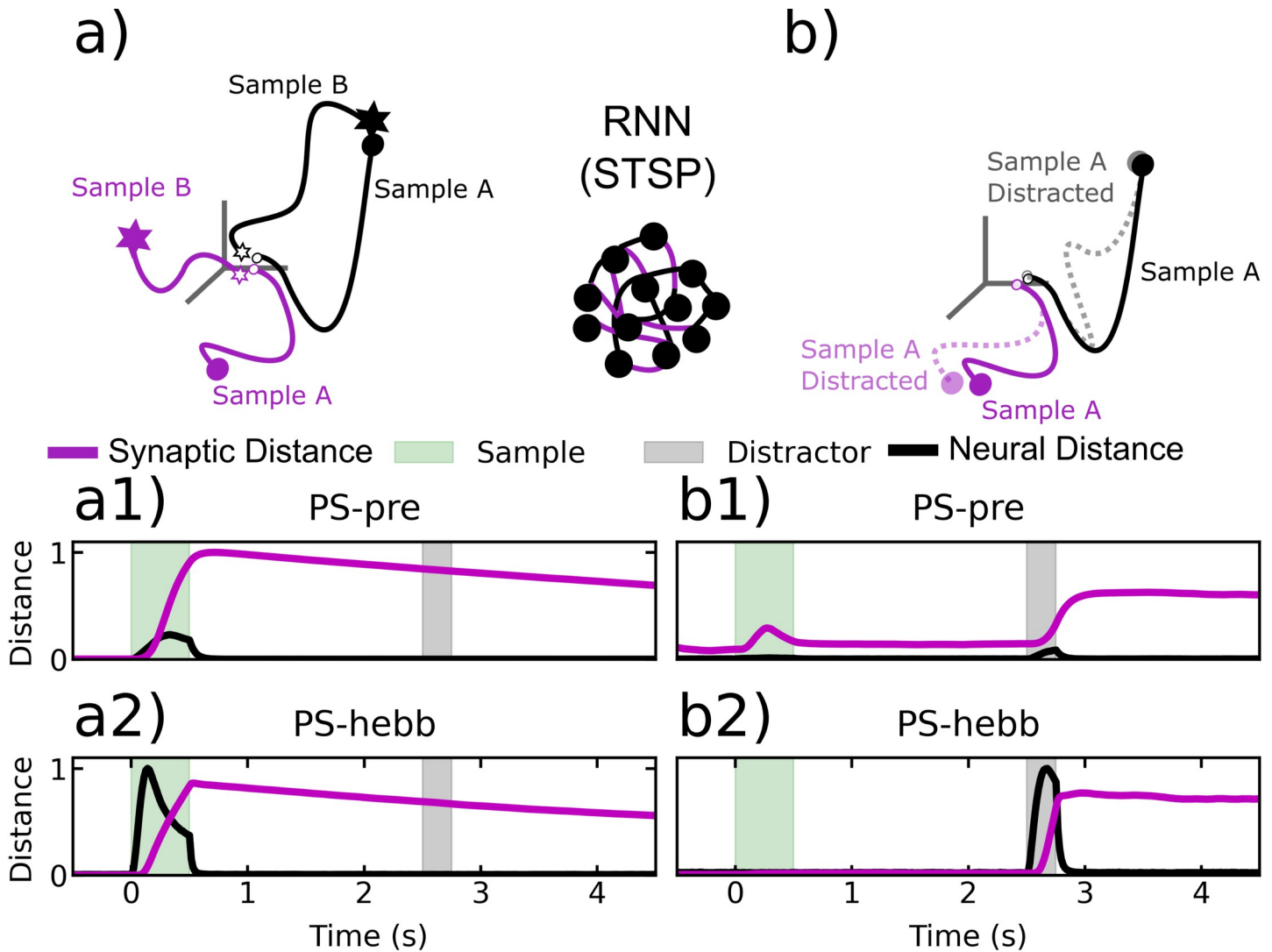
Comparison between trials with and without a distractor showed that the distractor temporarily knocked trajectories out of the attractor state. These trajectories quickly returned to the pre-distractor attractor state. This is how the fixed-synapse RNNs achieve robustness to the distractor. Once the neural trajectories were ‘captured’ by the appropriate attractor, they remained fixed there and can return there after being perturbed. In terms of neural activity this means that the spike rates stayed elevated across the delay period corresponding to the sample identity. In other words, these models achieve robust WM through persistent neural activity as a consequence of attractor dynamics, as observed in many prior models.



**Fig 7. Distances between neural trajectories within a sample condition and between sample conditions, for fixed synapse RNNs.** All models used were the most 'brain-like', as determined by the methodology in section "RNNs With STSP are More Brain-Like". a) Cartoon of trial-averaged RNN trajectories corresponding to two different sample conditions for the fixed-synapse RNNs. a1) Average pairwise distance between trajectories on different sample conditions, for the fixed synapse network with tanh activation (FS-tanh). a2) The same plot as in a1, but for FS-relu. b1) The average distance between distracted and undistracted trajectories. Average taken over all sample conditions. Results shown for FS-tanh. b2) Same results as in b1, but for FS-relu.

<https://doi.org/10.1371/journal.pcbi.1010776.g007>

The models with STSP did not achieve robust WM through persistent neural spiking, but instead relied on the STSP. We again examined trajectories using spikes rates from the STSP models (Fig 8). The neural distance between trajectories for different samples increased during sample presentation but then dropped back down near zero during the delay, especially longer delays. This mirrors the results from actual MUA activity in the PFC (see above). This stands in contrast to the fixed-synapse models (Fig 7), in which strong persistent spiking occurred. However, unlike the fixed-synapse models, in the STSP models we could measure trajectories and distances in synaptic state-space. This is analogous to using spike rates but instead we measure synaptic weights over time (as was done in Masse [36]). This revealed that trajectories for different samples in synaptic state-space remained elevated across the delay period (Fig 8), indicating sample information was being maintained by STSP even when spiking levels were low.



**Fig 8. Distances between neural and synaptic trajectories within a sample condition and between sample conditions, for plastic synapse RNNs.** Black lines correspond to neural trajectories, purple lines correspond to synaptic trajectories. a) Cartoon of trial-averaged neural and synaptic RNN trajectories corresponding to two different sample conditions for the plastic-synapse RNNs. a1) Average pairwise distance between neural and synaptic trajectories on different sample conditions, for PS-pre. a2) The same plot as in a1, but for PS-hebb. b1) The average distance between distracted and undistracted trajectories. Average taken over all sample conditions. Results shown for PS-pre. b2) Same results as in b1, but for PS-hebb.

<https://doi.org/10.1371/journal.pcbi.1010776.g008>

As we found in the experimentally recorded MUA, in the STSP models the spike rate trajectories between distractor and non-distractor trials increased during distractor presentation but then quickly decreased back to pre-distraction levels (Fig 8). By contrast, the distractor had a longer lasting effect in synaptic weights for the STSP models (Fig 8). These plots show the synaptic distance decreasing on a longer timescale than the neurons. Since the synaptic distance decreased on a timescale several times larger than the intrinsic time-constant of synapses, this indicates that the neurons and synapses interact in a way that increases their effective time constant [29]. This phenomenon is also found in simple linear systems, where a judicious choice of weight matrix (for example a marginally stable weight matrix) can lead to an increased effective time constant [29]. This could also potentially increase the susceptibility of the STSP networks to distractions. However, notably, we found that this was not the case. This

increased time constant did not affect the ability of the STSP networks to perform the task. The STSP networks performed at a high level (above 90% accuracy) on both distractor and non-distractor trials.

## Discussion

We found that while RNNs with and without STSP showed robustness against distractors, the RNNs with STSP were more brain-like. We also found that RNNs with STSP were more structurally robust than RNNs without STSP. STSP models showed graceful degradation to synaptic loss. In the STSP models, decodability of spike rates decreased during the memory delay (especially longer delays). Importantly, the WMs *could* be decoded from the synaptic weights. This was in contrast to the RNNs without STSP, where spike-rate decodability remained high over the entire memory delay. This high spike-rate decodability did not match observations of actual spiking recorded from the PFC. In sum, STSP can not only maintain information over gaps of no spiking, it also adds functional advantages. And, notably, adding STSP to RNNs makes them exhibit brain-like behavior.

Some forms of STSP worked better than others. We found that purely Hebbian STSP models were difficult to train, in agreement with previous studies [37]. By contrast, our anti-Hebbian STSP model, *PS-hebb*, (which reduced synaptic weights when spiking was too correlated) was successful and did not require any external weight clipping. We and others have found that anti-Hebbian STSP plays a crucial role in network stability and function [34,38]. And like others, we found that the synaptic weights in the *PS-pre* model had to be limited to a certain range to maintain stability and trainability [19,36]. This assumption is biologically plausible, as biological synapses are limited in the amount of resources that can be recruited by a presynaptic spike [19,39].

We found that the LSTM and GRU networks tended to be more robust to structural perturbations as well as process noise than the other networks. However, they were less brain-like than the STSP networks. This suggests that STSP may carry additional functional advantages in the brain beyond robustness, such as energy-efficiency (spikes are metabolically expensive). We also found that the STSP networks were more structurally robust than the fixed-synapse networks. This suggests that fixed-synapse networks require a “fine-tuning” of their synaptic weights, to appropriately mold their attractor landscape in a way that subserves the demands of the WM task. By contrast, we found that fixed-synapse networks tended to be more robust to process noise than the STSP networks. Given the superior brain-similarity of the STSP networks, this suggests that the brain may optimize for higher structural robustness as opposed to higher robustness to process noise. Indeed, the turnover rate of dendritic synapses in sensory and motor brain areas is as high as 40% every five days. Further exploring this link between STSP, structural robustness, and brain-similarity is an interesting direction for future studies.

The STSP models are in line with a variety of studies suggesting that cognition is more complex than steady attractor states. For example, sustained attention was long-thought to depend on attractor-like steady-state spiking. However, like studies of WM, this may have been an artifact of averaging spiking across trials [11,14]. Examination of sustained attention in real time (on single trials) has shown that, behaviorally, attention waxes and wanes rhythmically at theta [40]. There is a corresponding waxing and waning of spiking synchronized to LFP theta oscillations. Likewise, neural correlates of WM show sparse bursts of spiking linked to oscillatory dynamics when examined in real time [14,15].

Whether WM relies on persistent attractor dynamics (i.e., persistent spiking) alone or sparse spiking combined with STSP is of current interest [9–11,41,42]. A recent computational study has provided insight. Masse et al trained a WM model with STSP [36,43]. They found

that there was sparse spiking when WMs were being simply maintained and more spiking when WMs were being used and manipulated. This is consistent with other STSP-based WM models. For example, Lundqvist et al [14], found that spiking increases when WMs are being read out for use. This makes sense because the brain cannot read information from the synapses directly. Spikes are needed to “ping” the network to read information from the synaptic weights. Thus, models with STSP can work for both modes: Sparse spiking when WMs are being maintained and more spiking when WMs are being used.

It is also worth noting that studies that report higher levels of memory-delay spiking tend to be those that used spatial delayed response tasks. Spatial delayed response may involve more motor inhibition than WM per se [15]. The animal knows the forthcoming behavioral response and is inhibiting it while waiting for a “go” signal. It is like revving the engine while keeping your foot on the brake. By contrast, in WM tasks like delayed match-to-sample (i.e., the task used in this study), the behavioral response is not known until after the memory delay. They involve holding information for further computation, not inhibiting a behavioral response. Tasks like delayed match-to-sample are well known to produce lower levels of memory delay spiking [14,15,36,44] and thus require more than neural attractor states alone. That being said, replication of our findings in other animals and similar tasks is needed to strengthen our claims. Furthermore, while we limited ourselves to single-circuit WM mechanisms, it will be useful in future work to explore the potential role of other brain areas (such as hippocampus), other forms of synaptic plasticity, other computational mechanisms (such as rotational dynamics [45]) and perhaps even different cell types [46].

The structural robustness and provable stability [34] added by *PS-hebb*, the most brain-like network we trained, is critical not just to WM but to top-down control in general. A control system which is overly sensitive to perturbations (e.g distractors) would quickly propagate these disturbances to the rest of the brain, leading to errant behavior. Moreover, robustness and stability are intimately tied to modularity [47,48]. Stable networks can be linked with other stable networks in a way that preserves stability. Without that property, the whole network must be retrained when new modules are added to subserve complex or composite behavior. Our results suggest that STSP has a dual role. It can maintain information while simultaneously ensuring robustness in an energy-efficient manner [9]. We speculate that stability could be a key component to understanding the way specialized modules in the brain dynamically cooperate to form cohesive perceptions, plans, and actions.

## Methods

### Ethics statement

The non-human primate subject used in our experiments was a male rhesus macaque monkey (*Macaca mulatta*), aged 17. All procedures followed the guidelines of the Massachusetts Institute of Technology Committee on Animal Care and the National Institutes of Health (protocol approval number: 0322-029-25).

### Subject and task

As described in the main text, Utah arrays were implanted bilaterally in dorsolateral PFC (dlPFC) and ventrolateral PFC (vlPFC). The animal performed a distracted delayed match-to-sample task. Each trial began with visual fixation on the middle of the screen for 0.5s. Fixation was maintained throughout the trial until the behavioral response. The delay length was parametrically varied from 1–4 s in five logarithmic steps, randomly chosen each trial. At mid-delay a neutral distractor (1 of 2 possible objects never used as samples) was presented randomly on 50% of trials. During the multi-choice test the animal was allowed to freely saccade



between all objects on the screen. The final choice was indicated by fixating on it for at least one second.

## Data analysis methods

### Distances between neural trajectories

All neural spiking data was first smoothed using a 10ms Gaussian kernel. We refer to the smoothed spiking data as the firing rate. To quantify the difference in population activity between different conditions, we computed the distance between neural trajectories. We define a neural trajectory as a vector of firing rates for  $N$  recorded neurons evolving in time.

For the distractor vs non-distractor comparison, denote  $\bar{x}_t^{s,d,1} \in [0, \infty)^N$  the neural trajectory at time  $t$  on trial  $l \in \{1, \dots, L\}$  for stimulus identity  $s \in \{1, \dots, S\}$  and where  $d \in \{0, 1\}$  denotes the presence or absence of a distractor in the delay phase. We first computed trial averages as

$$\bar{x}_t^{s,d} = \frac{1}{L} \sum_{l=1}^L x_t^{s,d,l}$$

Then for each  $s$  we computed the distance between distracted and non-distracted trajectories, using the Euclidean norm, and took the average over stimuli:

$$d_t^{\text{dist}} = \frac{1}{S} \sum_{s=1}^S \|\bar{x}_t^{s,1} - \bar{x}_t^{s,0}\|^2$$

For the comparison between trajectories of different stimulus identities, we used only non-distractor trials ( $d = 0$ ). We first averaged over trials as above for each stimulus to get

$$\bar{x}_t^s = \frac{1}{L} \sum_{l=1}^L x_t^{s,0,l}$$

and calculated the mean distance between each pair of stimuli as:

$$d_t^{\text{stim}} = \frac{2}{S(S-1)} \sum_{1 \leq i < j \leq S} \|\bar{x}_t^i - \bar{x}_t^j\|^2$$

### Confidence intervals for trajectory distances

As our final measurement we calculated the mean over six experimental sessions for both  $d_t^{\text{dist}}$  and  $d_t^{\text{stim}}$ . To quantify uncertainty that captures both the session and trial variability, we used a non-parametric multi-level bootstrap. Briefly, for  $b \in \{1, \dots, B\}$  with  $B = 1000$ , we first sampled with replacement over sessions (that is, we randomly selected a session with replacement  $E$  times, where  $E$  is the total number sessions in the dataset—in our case  $E = 6$ ), and then for each resampled session, we resampled trials with replacement in the trial-average steps described above, yielding the bootstrap samples  $d_t^{\text{dist},b}$  and  $d_t^{\text{stim},b}$ . We took the 2.5<sup>th</sup> and 97.5<sup>th</sup> percentiles across  $b$  as our lower and upper values for the 95% confidence interval.

### Sample decoding

For sample decoding experiments, we first smoothed trials (as above, with 10ms Gaussian kernel), and then took average rates per neuron in 50ms time bins and attempted to decode sample type from the population activity at each time point. For our classifier, we used a linear

Support Vector Machine and standard regularization ( $C = 1$  in SciKit Learn). For each session we used a 10-fold cross-validation to calculate decoding accuracy. (That is, we held out 10% of trials as a test set, and trained on the remaining 90% of the data. We did this for 10 non-overlapping test sets, and then took the average test accuracy). We performed a decoding analysis on each session separately, and then took the average and standard error across sessions for our result.

### Change from baseline decoding

We smoothed trials and took average rates as in the sample decoding above. Then, for a given sample type, we took the population activity from 400-350ms before the stimulus (or distractor) and attempted to decode subsequent population activity against this baseline activity. We performed 10-fold cross-validation as above for calculating test accuracy, and calculated the average performance over samples for each session. Plots show the average and standard error across sessions.

### Sensitivity to smoothing and bin size

Throughout our data analysis we used a 10ms Gaussian kernel to produce firing rates, and in the decoding analysis, we used the mean firing rate in 50ms time bins. To ensure our trajectory distance results did not depend specifically on the choice of Gaussian kernel, we examined trajectory distances using 5, 10, and 15ms Gaussian kernels (S6 Fig). For the distance between trajectories with different sample identities, we took the mean distance within the sample presentation period (0, 0.5s) and the 100ms preceding test time ( $t_{test}-0.1s, t_{test}$ ), for each delay and each session, and then compared the sample distances and test distances. Similarly, for distances between distracted and non-distracted trajectories, we took the mean distance within the distractor presentation ( $t_{dist}, t_{dist}+0.25s$ ), and the ( $t_{test}-0.1s, t_{test}$ ), again for each delay and each session. To ensure our decoding results did not depend on the specific choice of Gaussian kernel and length of time bin, we again used 5, 10, and 15ms Gaussian kernels and additionally used 20ms, 50ms, and 100ms time bins (S7 Fig). We then looked at sample decoding accuracy in the non-distracted trials, for each combination of Gaussian kernel and time bin, again taking a mean for the sample presentation (0, 0.5s), and a mean for the pre-test period ( $t_{test}-0.1s, t_{test}$ ) for each delay and each session. The Wilcoxon signed-rank test was used for comparisons between time epochs (e.g. Sample vs Pre-Test), with p-value corrections for multiple comparisons made using the Benjamini-Hochberg procedure (with  $\alpha = 0.05$ ).

### Artificial neural network modelling

**Fixed-Synapse (FS).** For the fixed-synapse ANNs we used a ‘vanilla’ recurrent neural network. Each neuron had an activation  $x_i$  which we collected into a vector  $x \in R^n$  such that  $(x)_i = x_i$ . This vector was the *state* of the neural network at a given time  $t$ . The state evolves in time according to the following differential equation:

$$\tau \dot{x} = -x + \phi(Wx + I + b + Noise)$$

Here  $\tau$  is the time-constant of the network, which we set to 100 milliseconds. The function  $\phi(\cdot)$  is a nonlinearity (e.g. tanh, softplus, etc). The matrix  $W$  is the recurrent weight matrix of the network, and determines how the different neurons interact with one another. The vector  $I \in R^n$  is the exogenous input into the recurrent neural network, and can vary with  $t$ . For our purposes  $I$  will represent sensory information related to the task at hand. The vector  $b \in R^n$  is a fixed bias term. The noise will be defined shortly. To integrate the neural dynamics, we use the standard Euler approximation scheme with fixed step-size of  $dt = 15$  milliseconds. Defining the ratio of  $dt$

to  $\tau$  as  $\alpha \equiv \frac{dt}{\tau}$ , the approximated dynamics are:

$$x_{t+1} = (1 - \alpha)x_t + \alpha\phi(Wx_t + I_t + b + \text{Noise})$$

Here  $\text{Noise} = \sqrt{(2\alpha^{-1}\sigma_{rec}^2)}N(0, 1)$  where  $\sigma_{rec} = 0.05$  denotes the *process noise* and  $N(0,1)$  denotes a draw from the standard normal distribution. For the *FS-tanh* networks, we use  $\phi(\cdot) = \tanh(\cdot)$ . For the *FS-relu* networks, we use  $\phi(\cdot) = \text{ReLU}(\cdot) = \max(0, \cdot)$ . Finally, we used an affine readout for the network to produce the desired output for the task:

$$y = W_{out}x + c$$

With the exception of  $W$ , all the parameters of the fixed-synapse networks were initialized by drawing each element randomly from a gaussian distribution with mean zero and standard deviation  $\frac{1}{\sqrt{n}}$ . For the recurrent weight matrix  $W$  the elements were drawn from a gaussian distribution with mean zero, but standard deviation  $\frac{0.9}{\sqrt{n}}$ . We found that this initialization achieved faster training. At the beginning of each trial, the network state was initialized at the origin (i.e.  $x = 0$ ).

**Plastic synapse.** *PS-pre.* This ANN was an implementation of the model originally introduced by Mongillo et al [19] and then subsequently trained using deep learning techniques by Masse et al [36]. Many of the details are similar to or exactly the same as in Masse et al. We go through them here for completeness. As in the fixed-synapse model, we used Euler integration to approximate the dynamics:

$$x_{t+1} = (1 - \alpha)x_t + \alpha\phi(Wx_t u_t a_t + I_t + b)$$

Aside from the addition of two new variables  $u_t$  and  $a_t$  (which we'll define shortly) this is the same model as the fixed-synapse model. This network consisted of 80% excitatory neurons and 20% inhibitory neurons. To keep these populations separated throughout training, we decomposed  $W$  as  $W = W_{rec}^+ D$  where  $W_{rec}^+$  is a matrix with non-negative entries and  $D$  is a diagonal matrix whose  $i^{th}$  entry is either 1 or -1, depending on whether neuron  $i$  is an excitatory or inhibitory neuron. To each synapse we associated a value  $a$ , for the fraction of available neurotransmitter and  $u$ , the amount of utilized neurotransmitter. These two variables evolved according to:

$$\dot{a} = \frac{1 - a}{\tau_a} - u(t)a(t)x(t)dt$$

$$\dot{u} = \frac{U - u(t)}{\tau_u} + U(1 - u(t))x(t)dt$$

here  $x(t)$  is the presynaptic activity at time  $t$ ,  $\tau_a$  is the time constant of available neurotransmitter recovery,  $\tau_u$  is the time constant of neurotransmitter utilization recovery. Half the synapses in the network were facilitating, the other half was depressing. For the facilitating synapses,  $\tau_a = 200 \text{ ms}$ ,  $\tau_u = 1500 \text{ ms}$  and  $U = 0.15$ . For the depressing synapses,  $\tau_a = 1500 \text{ ms}$ ,  $\tau_u = 200 \text{ ms}$  and  $U = 0.45$ . Aside from  $W_{rec}^+$ , the parameter initialization for *PS-pre* was the same as the fixed synapse networks. For  $W$ , we found that initializing this matrix by drawing each element from a Gaussian with mean 0 and standard deviation  $1/N$ , where  $N$  is the number of hidden neurons, and then randomly setting 50% of the weights to zero, worked the best. To ensure that  $W_{rec}^+$  stayed non-negative throughout training, we always ran each element through a rectified linear function at the beginning of the trial. For each trial, the neural state vector  $x$  was initialized at the origin. The facilitation variables  $a$  were all initialized at  $a = 1$ . The utilization

variables were all initialized at  $u = 0$ . Throughout the trial, both facilitation and utilization were clamped between 0 and 1. The readout from the network was affine, as in the fixed-synapse networks.

*PS-hebb*. This excitatory anti-Hebbian / inhibitory Hebbian ANN was originally introduced in Kozachkov et al [34]. As before, we use a Euler integration scheme to approximate the dynamics:

$$x_{t+1} = (1 - \alpha)x_t + \alpha\phi(W_t x_t + I_t + b + \text{Noise})$$

$$W_{t+1} = (1 - \alpha\gamma)W_t - \alpha K \circ x x^T$$

Here  $\phi(\cdot) = \cdot$ . That is,  $\phi$  is just the identity function. The parameter  $\gamma$  controls the degree of weight decay, and is set to  $1/200\text{ms}^{-1}$  unless otherwise noted. The matrix  $K$  is positive and positive-definite, and the symbol  $\circ$  denotes element-wise matrix multiplication (i.e Hadamard product). During training we ensure positiveness and positive-definiteness of  $K$  by parameterizing it as:

$$K = B^T B + 10^{-2}O + 10^{-2}I$$

Where  $B_{ij} = C_{ij}^2$ ,  $O_{ij} = 1$  and  $I$  is the identity matrix. Since  $B^T B$  is positive semi-definite and non-negative, the matrix  $O$  ensures that  $K$  is strictly positive, while the identity matrix ensure that  $K$  is strictly positive-definite. The training is done on the matrix  $C$ . The output from the network is affine, as in all the other networks. All the parameters are initialized the same way as the fixed-synapse networks, except for  $C$ . The elements of this matrix are drawn from a uniform distribution between  $-0.5$  and  $0.5$ . At the beginning of each trial, the neural state vector  $x$  and the weights  $W$  are initialized at the origin.

For all models, we tested how the numerical integration step-size impacted performance. We used the hyperparameters corresponding to the most brain-like models. For each of the models, we tested a smaller step-size (10ms) and a larger step-size (30ms). All the models expect *PS-hebb* performed the same. *PS-hebb* was sensitive to the larger step-size, but insensitive to the smaller step-size. This makes sense because the stability guarantees for *PS-hebb* were proven in continuous-time. Smaller step-sizes provide more accurate approximations to the precise continuous-time solution.

We found that the initialization of the parameters was important for all models. For the fixed-synapse models, initializing the weight matrix with a slightly sub-critical matrix yielded a good tradeoff between initial stability and expressivity. For the *PS-pre* model, using the log-normal initialization provided by Masse worked well. When we experimented with different initializations, training either failed or was slower. The same was true of the *PS-hebb* model—different initialization schemes led to different training profiles in terms of model performance and training time.

**Task structure and training details for artificial neural networks.** The basic task structure for the ANNs was the same as for the non-human primate. Each stimuli was encoded in a one-hot vector. There was a total of 10 images (8 sample images, 2 distractor). The number of input neurons going into the ANNs was  $11 = 10 + 1$ : 10 input neurons for the stimuli and 1 for a fixation signal. These input stimuli were collected into a vector  $m_t$ , which changed with time throughout the trial. This vector was fed into the ANNs via a trainable input weight matrix  $I_t = W_{in} m_t$ . There were  $11 = 8 + 2 + 1$  output neurons leaving the ANNs, 8 for each sample image, 2 for the distractors and one as a fixation neuron. These neurons were collected into the vector  $y_{desired}$ , which defined the desired output of the ANNs.

The fixation is left on for 1000ms. At the end of the fixation period a randomly chosen sample image is presented for 500ms. Following the sample presentation period, the delay period begins. The delay length is randomly chosen from 1.0 s, 1.41s, 2.0 s, 2.83 s and 4s.

On 50% of the trials a distractor is presented in the middle of the delay for 250ms. Following the end of the delay period, the test period begins. During the test period the original sample image is presented as well as an off-target image for 500ms. At this point the RNN must do two things: stop fixating (i.e make the output fixation neuron output a value of '0') and indicate the correct choice by making the output neuron corresponding to the sample image hold a value of '1'.

We use a mean-squared loss to quantify network performance. For each trial we take the network output for a period of 500ms after the test period ends and compute the loss:

$$\mathcal{L} = \frac{1}{n_{out}T} \sum_t^{t+T} \|y - y_{desired}\|^2$$

We compute this loss for each trial in the batch, and then average over the batch to get a total loss per batch. Our training and test sets consisted of  $2^{14}$  trials each. We used a batch size of 256 to train and test. To train the networks we use the Adam optimizer with PyTorch standard parameters, a weight-decay value of  $10^{-4}$  for all models unless otherwise specified. For the FS models we use a learning rate of  $10^{-3}$ . For *PS-pre* we use a learning rate of 0.02. For *PS-hebb* we use a learning rate of 0.01.

As described in the main text, we concluded a hyperparameter search to assess the generality of our claims. For reasons we will describe shortly, the *PS-hebb* model was treated slightly differently than the other three. For *FS-relu*, *FS-tanh* and *PS-pre*, we looped over:

- Hidden size = (100, 129, 166, 215, 278, 359, 464, 599, 774, 1000)
- Activity regularization = (1e-03, 3.59e-03, 1.29e-02, 4.64e-02, 1.66e-01, 5.99e-01, 2.15e+00, 7.74, 2.78e+01, 1e+02)
- Parameter regularization (weight decay) = (1e-4, 1e-3, 1e-2)

For activity regularization, we computed the squared sum of all the neuron activations in a batch, and divided by the number of neurons and the number of time points. This quantity was then scaled by the regularization strength parameter which was varied during the hyperparameter sweep. This scaled quantity was then added to the overall loss function. For the *PS-hebb* model, we varied the hidden size by dividing the above hidden size values by 10 and finding the nearest integer. The reason we could not increase the size of the model to 1000 is due to GPU memory issues. The number of variables in *PS-hebb* model scales quadratically with the number of neurons (because each synapse is also a variable).

To test network performance, we compute the 'decision' of the network as the index of the output neuron which had the highest value. We compute the average accuracy of the network by the following procedure:

1. In the 500ms following the end of the test period, we compare the network output to the desired output for each time-point. We add up the number of times the network makes the right decision in this period, and then divide over the number of time-points to get a single-trial accuracy.
2. We compute this single-trial accuracy for all trials in the test set, and report the final network accuracy on the task as the average single-trial accuracy, taken over the whole test set.

## Supporting information

For all supplementary figures, green indicator box denotes sample period, gray distractor period, blue test period. Dotted lines indicate chance accuracy values

**S1 Fig. Decoding sample identity from neural spiking, Classifier was a linear support vector machine.** We evaluated performance on held-out data. Sample decoding accuracy drops to approximately chance during the delay.

(TIF)

**S2 Fig. Average distance between pairs of trajectories elicited by different samples.**

(TIF)

**S3 Fig. Decoding pre vs post sample presentation from neural spiking.** Spiking contains information about pre vs. post sample presentation through the delay. This indicates that neural spiking does not return to its pre-sample firing pattern following sample presentation.

(TIF)

**S4 Fig. Average distance between distracted and non-distracted trajectories corresponding to the same sample item.**

(TIF)

**S5 Fig. The effect of subsampling the RNNs on the brain-correlation index used.** The qualitative result stays the same (PS-hebb is more brain-like on average), but the particular correlation values change.

(TIF)

**S6 Fig. Effect of using different smoothing parameters on trajectory distances.** In all plots, each pair of dots correspond to averages during two time epochs for a particular session and delay time. In both cases the Pre-Test epoch used was 100ms preceding test time. All comparisons were statistically significant ( $p < 1e-3$ , Wilcoxon signed-rank test). a) Average distance between neural trajectories for different sample IDs for different smoothing parameters. Comparison is between Sample epoch (blue) and Pre-Test epoch (red). b) Average distance between neural trajectories for distracted and non-distracted trials for different smoothing parameters. Comparison is between Distractor epoch (blue) and Pre-Test epoch (red).

(TIF)

**S7 Fig. Average decoding accuracy on non-distractor trials for different smoothing parameters and different time bin sizes.** Comparison is between Sample epoch (blue) and Pre-Test epoch (red). In all plots, each pair of dots correspond to averages during two time epochs for a particular session and delay time. The Pre-Test epoch used was 100ms preceding test time. All comparisons were statistically significant ( $p < 1e-3$ , Wilcoxon signed-rank test).

(TIF)

**S8 Fig. The results of the hyperparameter sweep in the main text, including the LSTM and GRU networks.**

(TIF)

## Author Contributions

**Conceptualization:** Leo Kozachkov, John Tauber, Mikael Lundqvist, Scott L. Brincat, Jean-Jacques Slotine, Earl K. Miller.

**Data curation:** Mikael Lundqvist, Scott L. Brincat.



**Formal analysis:** Leo Kozachkov, John Tauber, Jean-Jacques Slotine.

**Funding acquisition:** Earl K. Miller.

**Investigation:** Leo Kozachkov.

**Methodology:** Earl K. Miller.

**Project administration:** Earl K. Miller.

**Resources:** Earl K. Miller.

**Validation:** John Tauber.

**Visualization:** Leo Kozachkov, John Tauber, Earl K. Miller.

**Writing – original draft:** Leo Kozachkov, John Tauber, Earl K. Miller.

**Writing – review & editing:** Leo Kozachkov, John Tauber, Mikael Lundqvist, Scott L. Brincat, Earl K. Miller.

## References

1. Baddeley A. Working memory. *Science*. 1992; 255: 556–559. <https://doi.org/10.1126/science.1736359> PMID: [1736359](https://pubmed.ncbi.nlm.nih.gov/1736359/)
2. Miller EK, Cohen JD. An integrative theory of prefrontal cortex function. *Annu Rev Neurosci*. 2001; 24: 167–202. <https://doi.org/10.1146/annurev.neuro.24.1.167> PMID: [11283309](https://pubmed.ncbi.nlm.nih.gov/11283309/)
3. Funahashi S, Bruce CJ, Goldman-Rakic PS. Mnemonic coding of visual space in the monkey's dorso-lateral prefrontal cortex. *J Neurophysiol*. 1989; 61: 331–349. <https://doi.org/10.1152/jn.1989.61.2.331> PMID: [2918358](https://pubmed.ncbi.nlm.nih.gov/2918358/)
4. Goldman-Rakic PS. Cellular basis of working memory. *Neuron*. 1995; 14: 477–485. [https://doi.org/10.1016/0896-6273\(95\)90304-6](https://doi.org/10.1016/0896-6273(95)90304-6) PMID: [7695894](https://pubmed.ncbi.nlm.nih.gov/7695894/)
5. Fuster JM, Alexander GE. Neuron activity related to short-term memory. *Science*. 1971; 173: 652–654. <https://doi.org/10.1126/science.173.3997.652> PMID: [4998337](https://pubmed.ncbi.nlm.nih.gov/4998337/)
6. Amit DJ, Brunel N. Model of global spontaneous activity and local structured activity during delay periods in the cerebral cortex. *Cereb Cortex N Y N 1991*. 1997; 7: 237–252. <https://doi.org/10.1093/cercor/7.3.237> PMID: [9143444](https://pubmed.ncbi.nlm.nih.gov/9143444/)
7. Hopfield JJ. Neural networks and physical systems with emergent collective computational abilities. *Proc Natl Acad Sci U S A*. 1982; 79: 2554–2558. <https://doi.org/10.1073/pnas.79.8.2554> PMID: [6953413](https://pubmed.ncbi.nlm.nih.gov/6953413/)
8. Kamiński J, Rutishauser U. Between persistently active and activity-silent frameworks: novel vistas on the cellular basis of working memory. *Ann N Y Acad Sci*. 2020; 1464: 64–75. <https://doi.org/10.1111/nyas.14213> PMID: [31407811](https://pubmed.ncbi.nlm.nih.gov/31407811/)
9. Stokes MG. 'Activity-silent' working memory in prefrontal cortex: a dynamic coding framework. *Trends Cogn Sci*. 2015; 19: 394–405. <https://doi.org/10.1016/j.tics.2015.05.004> PMID: [26051384](https://pubmed.ncbi.nlm.nih.gov/26051384/)
10. Barbosa J, Stein H, Martinez RL, Galan-Gadea A, Li S, Dalmau J, et al. Interplay between persistent activity and activity-silent dynamics in the prefrontal cortex underlies serial biases in working memory. *Nat Neurosci*. 2020; 23: 1016–1024. <https://doi.org/10.1038/s41593-020-0644-4> PMID: [32572236](https://pubmed.ncbi.nlm.nih.gov/32572236/)
11. Lundqvist M, Herman P, Miller EK. Working Memory: Delay Activity, Yes! Persistent Activity? Maybe Not. *J Neurosci*. 2018; 38: 7013–7019. <https://doi.org/10.1523/JNEUROSCI.2485-17.2018> PMID: [30089640](https://pubmed.ncbi.nlm.nih.gov/30089640/)
12. Pasternak T, Greenlee MW. Working memory in primate sensory systems. *Nat Rev Neurosci*. 2005; 6: 97–107. <https://doi.org/10.1038/nrn1603> PMID: [15654324](https://pubmed.ncbi.nlm.nih.gov/15654324/)
13. Cromer JA, Roy JE, Miller EK. Representation of Multiple, Independent Categories in the Primate Prefrontal Cortex. *Neuron*. 2010; 66: 796–807. <https://doi.org/10.1016/j.neuron.2010.05.005> PMID: [20547135](https://pubmed.ncbi.nlm.nih.gov/20547135/)
14. Lundqvist M, Rose J, Herman P, Brincat SL, Buschman TJ, Miller EK. Gamma and Beta Bursts Underlie Working Memory. *Neuron*. 2016; 90: 152–164. <https://doi.org/10.1016/j.neuron.2016.02.028> PMID: [26996084](https://pubmed.ncbi.nlm.nih.gov/26996084/)
15. Miller EK, Lundqvist M, Bastos AM. Working Memory 2.0. *Neuron*. 2018; 100: 463–475. <https://doi.org/10.1016/j.neuron.2018.09.023> PMID: [30359609](https://pubmed.ncbi.nlm.nih.gov/30359609/)

16. Constantinidis C, Funahashi S, Lee D, Murray JD, Qi X-L, Wang M, et al. Persistent Spiking Activity Underlies Working Memory. *J Neurosci*. 2018; 38: 7020–7028. <https://doi.org/10.1523/JNEUROSCI.2486-17.2018> PMID: 30089641
17. Wong K-F, Wang X-J. A Recurrent Network Mechanism of Time Integration in Perceptual Decisions. *J Neurosci*. 2006; 26: 1314–1328. <https://doi.org/10.1523/JNEUROSCI.3733-05.2006> PMID: 16436619
18. Lundqvist M, Herman P, Lansner A. Theta and Gamma Power Increases and Alpha/Beta Power Decreases with Memory Load in an Attractor Network Model. *J Cogn Neurosci*. 2011; 23: 3008–3020. [https://doi.org/10.1162/jocn\\_a\\_00029](https://doi.org/10.1162/jocn_a_00029) PMID: 21452933
19. Mongillo G, Barak O, Tsodyks M. Synaptic Theory of Working Memory. *Science*. 2008; 319: 1543–1546. <https://doi.org/10.1126/science.1150769> PMID: 18339943
20. Sandberg A, Tegnér J, Lansner A. A working memory model based on fast Hebbian learning. *Netw Bristol Engl*. 2003; 14: 789–802. PMID: 14653503
21. Knoblauch A, Palm G, Sommer FT. Memory Capacities for Synaptic and Structural Plasticity. *Neural Comput*. 2010; 22: 289–341. <https://doi.org/10.1162/neco.2009.08-07-588> PMID: 19925281
22. Seeholzer A, Deger M, Gerstner W. Stability of working memory in continuous attractor networks under the control of short-term plasticity. *PLOS Comput Biol*. 2019; 15: e1006928. <https://doi.org/10.1371/journal.pcbi.1006928> PMID: 31002672
23. Taher H, Torcini A, Olmi S. Exact neural mass model for synaptic-based working memory. *PLOS Comput Biol*. 2020; 16: e1008533. <https://doi.org/10.1371/journal.pcbi.1008533> PMID: 33320855
24. Bouchacourt F, Buschman TJ. A Flexible Model of Working Memory. *Neuron*. 2019; 103: 147–160.e8. <https://doi.org/10.1016/j.neuron.2019.04.020> PMID: 31103359
25. Durstewitz D, Seamans JK, Sejnowski TJ. Neurocomputational models of working memory. *Nat Neurosci*. 2000; 3: 1184–1191. <https://doi.org/10.1038/81460> PMID: 11127836
26. Jacob SN, Nieder A. Complementary Roles for Primate Frontal and Parietal Cortex in Guarding Working Memory from Distractor Stimuli. *Neuron*. 2014; 83: 226–237. <https://doi.org/10.1016/j.neuron.2014.05.009> PMID: 24991963
27. Barak O, Tsodyks M. Working models of working memory. *Curr Opin Neurobiol*. 2014; 25: 20–24. <https://doi.org/10.1016/j.conb.2013.10.008> PMID: 24709596
28. Murray JD, Bernacchia A, Freedman DJ, Romo R, Wallis JD, Cai X, et al. A hierarchy of intrinsic time-scales across primate cortex. *Nat Neurosci*. 2014; 17: 1661–1663. <https://doi.org/10.1038/nn.3862> PMID: 25383900
29. Dayan P, Abbott LF. *Theoretical Neuroscience: Computational and Mathematical Modeling of Neural Systems*. MIT Press; 2005.
30. Maheswaranathan N, Williams A, Golub M, Ganguli S, Sussillo D. Universality and individuality in neural dynamics across large populations of recurrent networks. *Advances in Neural Information Processing Systems*. Curran Associates, Inc.; 2019. Available: <https://proceedings.neurips.cc/paper/2019/hash/5f5d472067f77b5c88f69f1bcfda1e08-Abstract.html>
31. Hochreiter S, Schmidhuber J. Long Short-Term Memory. *Neural Comput*. 1997; 9: 1735–1780. <https://doi.org/10.1162/neco.1997.9.8.1735> PMID: 9377276
32. Cho K, van Merriënboer B, Bahdanau D, Bengio Y. On the Properties of Neural Machine Translation: Encoder-Decoder Approaches. *arXiv*; 2014. <https://doi.org/10.48550/arXiv.1409.1259>
33. O'Reilly RC, Frank MJ. Making Working Memory Work: A Computational Model of Learning in the Prefrontal Cortex and Basal Ganglia. *Neural Comput*. 2006; 18: 283–328. <https://doi.org/10.1162/089976606775093909> PMID: 16378516
34. Kozachkov L, Lundqvist M, Slotine J-J, Miller EK. Achieving stable dynamics in neural circuits. *PLOS Comput Biol*. 2020; 16: e1007659. <https://doi.org/10.1371/journal.pcbi.1007659> PMID: 32764745
35. Paszke A, Gross S, Massa F, Lerer A, Bradbury J, Chanan G, et al. PyTorch: An Imperative Style, High-Performance Deep Learning Library. *Advances in Neural Information Processing Systems*. Curran Associates, Inc.; 2019. Available: <https://proceedings.neurips.cc/paper/2019/hash/bdbca288fee7f92f2bfa9f7012727740-Abstract.html>
36. Masse NY, Yang GR, Song HF, Wang X-J, Freedman DJ. Circuit mechanisms for the maintenance and manipulation of information in working memory. *Nat Neurosci*. 2019; 22: 1159–1167. <https://doi.org/10.1038/s41593-019-0414-3> PMID: 31182866
37. Orhan AE, Ma WJ. A diverse range of factors affect the nature of neural representations underlying short-term memory. *Nat Neurosci*. 2019; 22: 275–283. <https://doi.org/10.1038/s41593-018-0314-y> PMID: 30664767

38. Tyulmankov D, Yang GR, Abbott LF. Meta-learning synaptic plasticity and memory addressing for continual familiarity detection. *Neuron*. 2021 [cited 24 Dec 2021]. <https://doi.org/10.1016/j.neuron.2021.11.009> PMID: 34861149
39. Tsodyks MV, Markram H. The neural code between neocortical pyramidal neurons depends on neurotransmitter release probability. *Proc Natl Acad Sci*. 1997; 94: 719–723. <https://doi.org/10.1073/pnas.94.2.719> PMID: 9012851
40. Helfrich RF, Fiebelkorn IC, Szczepanski SM, Lin JJ, Parvizi J, Knight RT, et al. Neural mechanisms of sustained attention are rhythmic. *Neuron*. 2018; 99: 854–865. <https://doi.org/10.1016/j.neuron.2018.07.032> PMID: 30138591
41. Parthasarathy A, Herikstad R, Bong JH, Medina FS, Libedinsky C, Yen S-C. Mixed selectivity morphs population codes in prefrontal cortex. *Nat Neurosci*. 2017; 20: 1770–1779. <https://doi.org/10.1038/s41593-017-0003-2> PMID: 29184197
42. Parthasarathy A, Tang C, Herikstad R, Cheong LF, Yen S-C, Libedinsky C. Time-invariant working memory representations in the presence of code-morphing in the lateral prefrontal cortex. *Nat Commun*. 2019; 10: 4995. <https://doi.org/10.1038/s41467-019-12841-y> PMID: 31676790
43. Masse NY, Rosen MC, Freedman DJ. Reevaluating the Role of Persistent Neural Activity in Short-Term Memory. *Trends Cogn Sci*. 2020; 24: 242–258. <https://doi.org/10.1016/j.tics.2019.12.014> PMID: 32007384
44. Miller EK, Erickson CA, Desimone R. Neural Mechanisms of Visual Working Memory in Prefrontal Cortex of the Macaque. *J Neurosci*. 1996; 16: 5154–5167. <https://doi.org/10.1523/JNEUROSCI.16-16-05154.1996> PMID: 8756444
45. Libby A, Buschman TJ. Rotational dynamics reduce interference between sensory and memory representations. *Nat Neurosci*. 2021; 24: 715–726. <https://doi.org/10.1038/s41593-021-00821-9> PMID: 33821001
46. Frontiers | Modeling Working Memory in a Spiking Neuron Network Accompanied by Astrocytes. [cited 12 Oct 2022]. Available: <https://www.frontiersin.org/articles/10.3389/fncel.2021.631485/full>
47. Slotine JE. Modular Stability Tools for Distributed Computation and Control. 2002.
48. Kozachkov L, Ennis M, Slotine J-J. Recursive Construction of Stable Assemblies of Recurrent Neural Networks. *ArXiv210608928 Cs Math Q-Bio*. 2021 [cited 24 Dec 2021]. Available: <http://arxiv.org/abs/2106.08928>