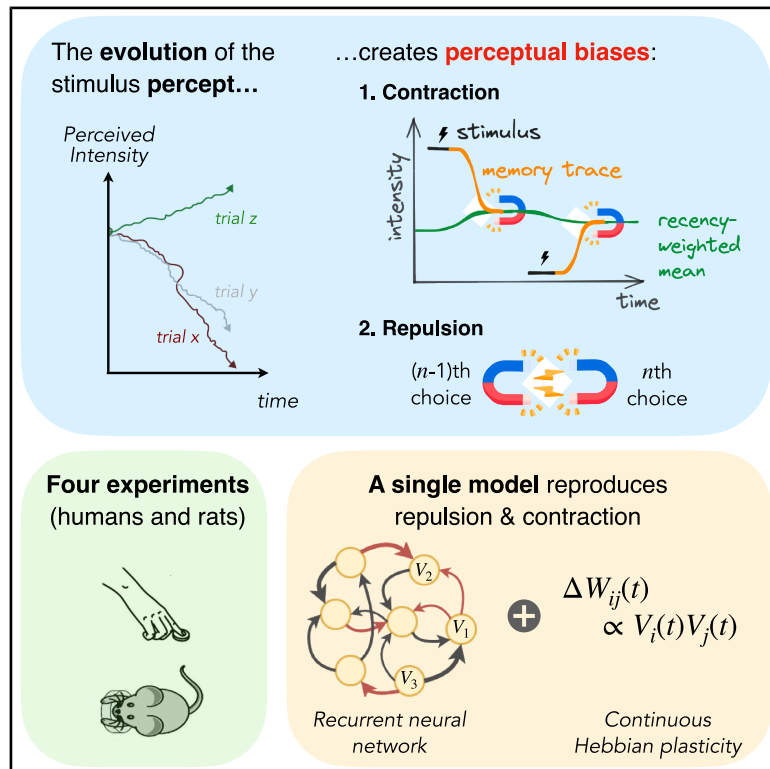


Diverse perceptual biases emerge from Hebbian plasticity in a recurrent neural network model

Graphical abstract



Authors

Francesca Schönsberg, Davide Giana, Yukti Chopra, Mathew E. Diamond, Sebastian Goldt

Correspondence

francesca.schonsberg@phys.ens.fr (F.S.), diamond@sissa.it (M.E.D.), sgoldt@sissa.it (S.G.)

In brief

Schönsberg et al. show that a simple recurrent neural network with continuous Hebbian plasticity reproduces classic perceptual biases like contraction and repulsion. Tested across three different tasks, the model matches experimental data without fine-tuning, suggesting that distinct biases can arise from a unified local learning mechanism in the brain.

Highlights

- Perceptual biases offer a glimpse into how the brain processes sensory stimuli
- We propose a single recurrent network model for diverse perceptual biases
- Contractive and repulsive biases emerge from ongoing Hebbian plasticity
- The model accurately predicts biases in four datasets covering three tasks



Article

Diverse perceptual biases emerge from Hebbian plasticity in a recurrent neural network model

Francesca Schönsberg,^{1,2,*} Davide Giana,³ Yukti Chopra,³ Mathew E. Diamond,^{3,4,*} and Sebastian Goldt^{2,4,5,*}¹Laboratory of Physics of the Ecole Normale Supérieure, PSL and CNRS UMR8023, Sorbonne Université, Paris, France²Department of Physics, International School for Advanced Studies (SISSA), Trieste, Italy³SENSEx Laboratory, International School for Advanced Studies (SISSA), Trieste, Italy⁴These authors contributed equally⁵Lead contact

*Correspondence: francesca.schonsberg@phys.ens.fr (F.S.), diamond@sisssa.it (M.E.D.), sgoldt@sisssa.it (S.G.)

<https://doi.org/10.1016/j.neuron.2025.09.037>

SUMMARY

Perceptual biases offer a glimpse into how the brain processes sensory stimuli. While psychophysics has uncovered systematic biases such as contraction (stored information shifts toward a central tendency) and repulsion (the current percept shifts away from recent percepts), a unifying neural network model for how such seemingly distinct biases emerge from learning is lacking. Here, we show that both contractive and repulsive biases emerge from continuous Hebbian plasticity in a single recurrent neural network.

We test the model on four datasets covering two sensory modalities in two working memory tasks, a reference memory task, and a novel “one-back task” designed to test the robustness of the model.

We find excellent agreement between model predictions and experimental data without fine-tuning the model to any particular paradigm.

These results show that apparently contradictory perceptual biases can emerge from a simple local learning rule in a single recurrent region of the brain.

INTRODUCTION

The seemingly simple task of categorizing the intensity of a stimulus is prone to several perceptual biases. For example, in a tactile intensity working memory task, where the subject is presented with two sequential stimuli on each trial and is required to report whether the second stimulus was weaker or stronger than the first (Figure 1A), both human and rodent subjects display a *contraction bias*: they tend to overestimate the strength of the first stimulus if it lies below the weighted average of past stimuli, and they underestimate the first stimulus if it lies above the average.^{1–4} A similar contraction occurs in tasks where subjects are asked to report the size or intensity of a stimulus after some time has elapsed since exposure to the stimulus.^{5,6} By contrast, in intensity reference memory tasks,^{7–10} where the subject reports whether one stimulus per trial was perceived as strong or weak, a *repulsive bias* emerges: subjects are more likely to characterize a stimulus as strong if the previous trial's stimulus was weak, and vice versa.¹⁰

These biases are of interest because they offer a window into the computational processes underlying perceptual judgments. On the experimental side, a great number of psychophysical studies, comprising different sensory modalities, species, and details of experimental design, have provided detailed characterization of both contraction and repulsion.^{1,3,5,7–11} On the

modeling side, comprehensive phenomenological models have been built to account for these biases.^{1,10,12–16} Nevertheless, the neural mechanisms at work remain poorly understood. While the idea that neural circuits could perform working memory tasks through a quasi-continuous line of fixed points dates back to the 1990s,^{17,18} Boboeva et al.¹⁹ showed only recently that a neural network model can reproduce the specific perceptual biases observed by Akrami et al.³ However, existing theoretical models have static ad hoc connectivity and are tuned to reproduce a single specific bias. Further, they overlook learning—the boundary must be learned in reference memory categorization, as evidenced by the notable performance disparities of trained subjects between the onset and the end of experimental sessions¹⁰; current models have failed to capture how perceptual biases arise in parallel with short-timescale learning. Finally, while recent phenomenological models point to the existence of general principles underlying diverse perceptual biases,¹⁶ there exists no single neural network model that, without reconfiguration, can account for multiple biases expressed across different experimental paradigms.

Here, we show that a single recurrent neural network (RNN) with ongoing Hebbian plasticity learns stimulus representations that reflect the perceptual biases of humans and rodents. We first show how the model reproduces contraction and repulsive biases, as well as the history dependence of choices, as



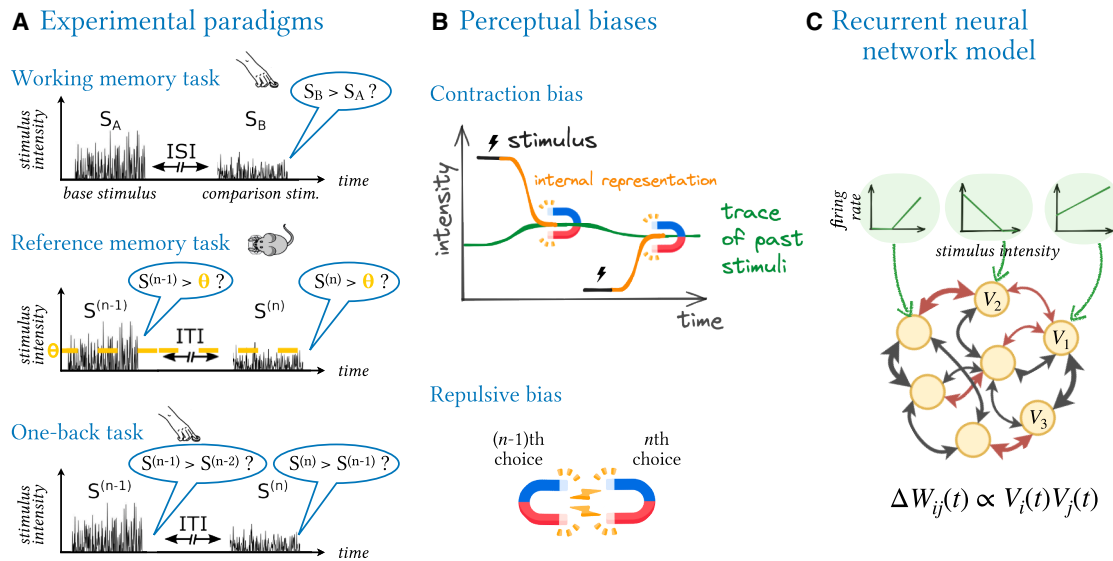


Figure 1. A single neural network model explains two perceptual biases in three perceptual memory paradigms

(A) We consider three experimental paradigms: a *working memory* task in humans, where subjects report whether the second of a pair of sequential stimuli of strengths S_A, S_B , separated by a variable inter-stimulus interval (ISI), was stronger than the first; a *reference memory* task where rats compare the strength of one stimulus per trial to a fixed category boundary (orange); and a novel *one-back task* designed to test the generality of our model. Here, human subjects compare the strength of one stimulus per trial to the strength of the stimulus in the preceding trial. In all three paradigms, trials are separated by a variable inter-trial interval (ITI), and subjects receive rewards after correct choices.

(B) The performance of humans and rodents in these paradigms reveals two perceptual biases: the *contraction* of stimulus representations to the trace of past stimuli and the *repulsion* between successive judgments.

(C) We reproduce both perceptual biases in a single neural network model, where a fully connected recurrent neural network (RNN) is driven by external inputs (green) encoding stimulus intensity, while recurrent connections are continuously reshaped by Hebbian plasticity and can be positive or negative (red/black).

measured in a novel tactile working memory experiment that we perform with human subjects. In addition, we show that our model also reproduces the biases observed in two previously published datasets, i.e., the auditory working memory task of Akrami et al.³ and the tactile reference memory task of Hachen et al.¹⁰ Our analysis shows how both perceptual biases emerge from a plastic attractor driven by recurrent dynamics and Hebbian plasticity. Finally, we design a new experimental paradigm that tests the generality of the model. In this one-back task [Figure 1A](#), participants have to compare the strength of the stimulus presented on each trial to the stimulus of the previous trial, making each stimulus serve both as comparison and, successively, as base. We find a novel, bimodal contraction bias in the performance of human subjects and reproduce it in the RNN model. We stress that we do not fine-tune the model to any individual paradigm, nor do we fit the RNN to experimental data or use gradient descent for learning. Instead, we fix hyper-parameters across experiments and use only local Hebbian plasticity to let the weights evolve continuously.

RESULTS

As an overview, we first illustrate the key experimental paradigms, the perceptual biases that emerge in such paradigms, and the schematic architecture of the neural network model ([Figure 1](#)).

An RNN model with ongoing Hebbian plasticity Network architecture

We considered an RNN model, taking inspiration from neurobiological evidence indicating that perceptual biases may occur in a recurrent subregion of the motor cortex (see [STAR Methods](#) for details). The network is composed of N threshold linear units, which model the neurons or the ensemble activity of a group of them. The activity state of the network $\vec{V}(t)$ represents the firing rate of all N units at each time t . Although the network state evolves continuously through the synchronous update of all units, for clarity, we present the update for each unit V_i individually, which obeys the equation^{20–22}:

$$V_i(t+1) = (1 - \gamma)V_i(t) + \gamma g \left[\frac{1}{N} \sum_{j=0}^N W_{ij}(t)V_j(t) - T(\langle \vec{V}(t) \rangle) + \eta \xi_i(t) \right]^+, \quad (\text{Equation 1})$$

where γ sets the timescale of the neuronal dynamics and $[z]^+ = \max(0, z)$. The neural input field, within the square brackets, is based on Treves²³ and Schönsberg et al.²⁴ and determines the activity of the neurons through three components: the inputs from the other neurons weighed by the recurrent connectivity $W_{ij}(t)$ (which represents the strength of the synapse from neuron j to neuron i), a time-dependent threshold $T(\langle \vec{V}(t) \rangle)$, and an external input $\xi_i(t)$, described below. The threshold does not

represent any specific biophysical quantity directly; instead, it serves to maintain a tendency toward a desired average activity value, preventing activities from diverging toward zero or infinity. Its value depends on the average activity across all units,

$$T(\langle \bar{V}(t) \rangle) \equiv 4\kappa \left(\frac{1}{N} \sum_i V_i(t) - a \right)^3 \text{ as in Schönsberg et al.,}^{24}$$

where a and κ are parameters setting the desired average and the contribution of the threshold, respectively (see STAR Methods for their values), while the cubic function is a simple supra linear function that preserves the sign of the deviation.

Synaptic plasticity

For simplicity, we assumed that the network is dense, with bidirectional connections $W_{ij}(t)$ between all pairs of neurons, defined in a connectivity matrix \bar{W} . All the connections evolve in parallel to the neural dynamics. Specifically, starting from random initial weights, the weights follow the Hebbian plasticity rule,²⁵ here reported again for one single connection:

$$W_{ij}(t+1) = \left(1 - \lambda - V_i(t)^2\right) W_{ij}(t) + (V_i(t) - a)(V_j(t) - a). \quad (\text{Equation 2})$$

The parameters λ and $V_i(t)^2$ are simple normalization terms that ensure the stable operation of a basic Hebbian learning rule, with the latter inspired by Oja's normalization.²⁶ The activity values in the second term, instead, are thresholded to yield both positive and negative weights. The key point is that the weights follow a biologically plausible learning rule based on Hebb's principle of synapse reinforcement, rather than being optimized using a gradient-descent-based machine learning algorithm.²⁷ The key distinctions between these two approaches lie in locality, continuity, and supervision. The Hebbian learning rule is both unsupervised and local, with changes in synaptic strength occurring solely due to the relationship of activity at presynaptic and postsynaptic neurons. These changes are gradual and continuous, driven by the network's ongoing activity. In contrast, the optimization-based approaches that our model avoids are relying on substantial amounts of information to train weights globally, requiring supervision and multiple iterations. While the Hebbian dynamics of Equation 2 may change the sign of some weights in an apparent violation of Dale's principle, we observed that this occurs infrequently (see section Synaptic Weight Changes and Biological Plausibility for details). We also note that recent experimental observations have opened the possibility of a switch in the excitatory/inhibitory effect of a neuron on its postsynaptic target.²⁸

We fix the hyper-parameters $\gamma, \lambda, \eta, g, \kappa, a$ to values that put the network in a regime where the recurrent dynamics and the learning occur without divergences (see STAR Methods for details). We stress that the model will be tested against the experimental results of all four datasets with the same set of parameters.

Stimulating the RNN

At each time step of the neural dynamics of Equation 1, the external stimulus delivered to the units $\vec{\xi}(t)$ can either be $\vec{0}$ (if no stimulus is applied) or a pattern that represents the intensity of the vibrational or auditory stimulus from the experiment. We designed a set of input patterns $\vec{\xi}^\mu$ for representing stimulus in-

intensities between S^{MIN} and S^{MAX} where the index $\mu = 1, \dots, P$ runs over the different discrete values of the stimulus intensity. In a scheme that reflects key coding properties of the rat's primary vibrissal somatosensory cortex (vS1),^{29,30} the input ξ_i^μ to the i th RNN unit for a stimulus of strength S^μ is a threshold linear function of the stimulus strength with a random positive or negative slope, as shown in Figure 1C (green inputs) and STAR Methods. This choice reflects the presence of neurons in vS1 whose tuning curves are either positively correlated or anti-correlated with stimulus strength. These conditions are similar to the binary morphing patterns of Blumenfeld et al.³¹ in that our inputs $\vec{\xi}^\mu$, even though continuous rather than binary, also display gradual morphing: the correlation between two patterns $\vec{\xi}^\mu$ and $\vec{\xi}^\nu$ decays with the absolute difference in the corresponding stimulus intensities $\Delta S = |S^\mu - S^\nu|$. This morphing structure will be key in shaping the dynamics of the RNN, as we discuss below. The precise stimulation protocols, i.e., the duration of each stimulation and each inter-stimulus interval (ISI), for all the tasks are given in appendix input patterns.

Reading out the intensity encoded by the RNN and detecting the choice

While the readout of the stimulus intensity represented by the RNN is independent of the task, the way we model the choice depends necessarily on the nature of the experiment. Inasmuch as the decision-making derives in humans by verbal instruction and in rats by an extensive training regime, it is reasonable to assume that learning of the behavioral context cannot be captured by model parameters in the current implementation. Specifically, we read out the intensity encoded by the RNN state $\vec{V}(t)$ by computing the normalized cosine similarity between $\vec{V}(t)$ and all the P input patterns $\vec{\xi}^\mu$. This "overlap profile" has P values that, owing to the correlation structure of the input patterns, follow a smooth, bell-shaped distribution with a single maximum. The stimulus intensity corresponding to the pattern $\vec{\xi}$ with the maximal overlap is referred to as the network's internal representation. The choice is instead determined based on the behavioral task. If the task requires comparing the current stimulus to a previous stimulus, as in the working memory and one-back experiments, we extract the choice at the second stimulus onset by comparing the second stimulus to the network's internal representation, which will reflect the previous stimulus. If the task instead requires categorizing the intensity of a single stimulus, as in the reference memory experiments, the choice is extracted by considering how the internal representation evolves after the stimulus is removed. If the internal representation evolves toward higher stimulus intensities, the stimulus is considered weak, and vice versa.

Human participants and the RNN model display contraction bias in a working memory task

We first applied the model to a human working memory experiment (Figure 1A). Participants were presented with pairs of vibratory stimuli of strengths (S_A, S_B) and tasked with reporting which of the two stimuli was stronger, earning a reward for correct choices (see STAR Methods for details). This paradigm has been used in several forms to explore working memory.^{2,3,18}

We conducted the experiment with 16 human subjects, each completing approximately 1,000 trials. The performance of

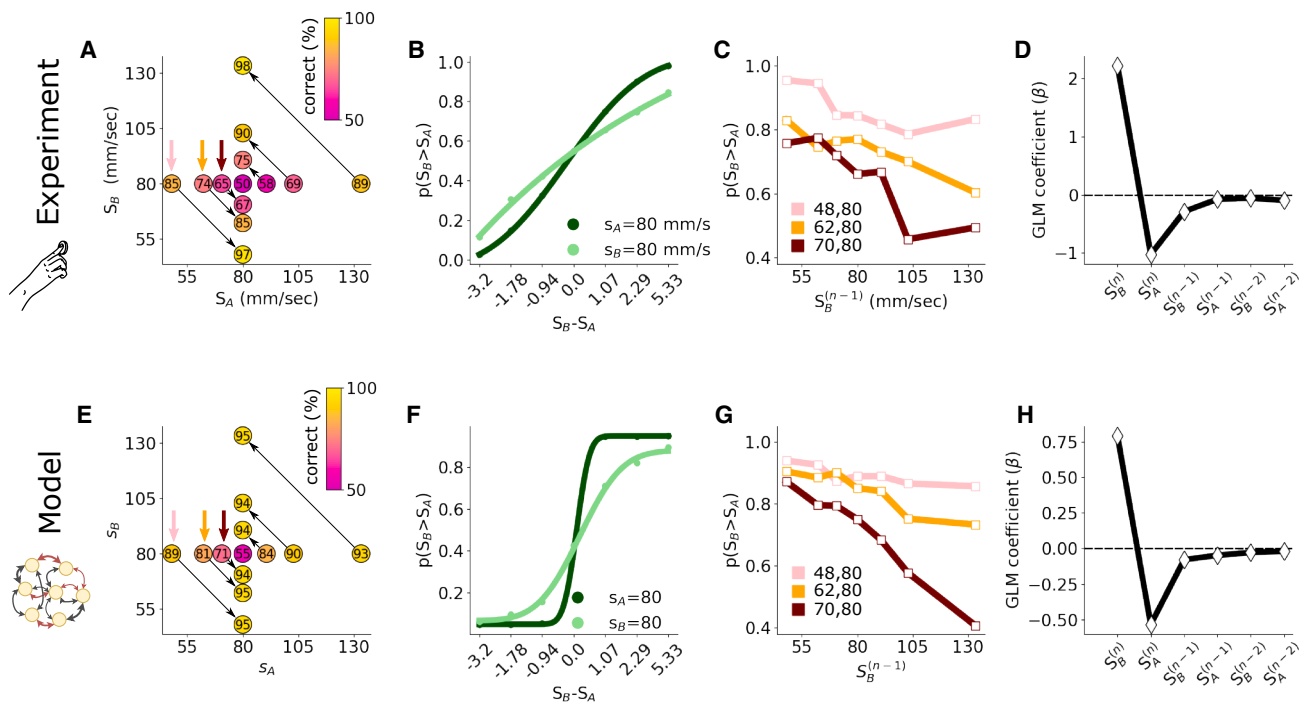


Figure 2. Contraction bias and history dependence of choices in a working memory task

(A) Accuracy of human participants in comparing the intensities of two vibrations (S_A, S_B). Arrows explained in main text. (B) Psychometric curves based on the stimulus difference if either the first or the second stimulus is 80 mm/sec (light and dark green, respectively). Points correspond to experimental data, solid lines to a sigmoidal fit (see STAR Methods). (C) Probability that the second stimulus is (correctly) classified as stronger as a function of the strength of the second stimulus S_B in the preceding trial. This plot is shown for the three stimulus pairs highlighted by colored arrows in (A). (D) Coefficient values (β) of a GLM fit of the choice at the n th trial as a function of the most recent six perceived stimuli. (E–H) The same results are reproduced in the RNN model under conditions mimicking those of (A)–(D).

human participants (Figure 2A) shows intriguing asymmetries. For example, the pair (105, 80) yielded 21% lower performance (69% vs. 90%) than the pair (80, 105), even if the difference between stimuli is of the same magnitude (25 mm/s). This asymmetry can be interpreted as an outcome of *contraction bias*^{3,5} as follows: during the ISI, the memory of the intensity of the base stimulus S_A moves toward an average of all previous stimuli, weighted by their recency through an unknown function. Thus, for the first pair, the memory of S_A , ideally fixed at 105, moves closer to the stimulus to which it must be compared, i.e., S_B (80), making the task harder and leading to errors. For the second pair, the memory of S_A moves toward the recency-weighted average, which is itself in the neighborhood of 80mm/sec. Because the base stimulus is in the vicinity of the recency-weighted average, little contraction occurs, leading to fewer errors. Indeed, the S_A memory may be “stabilized.” The thin black arrows in Figure 2A point toward the stimulus pair that is easier to discriminate among pairs with the same difference in stimulus strength. This bias also manifests in the average psychometric curves across subjects, after separating trials according to whether $S_A = 80$ mm/s or $S_B = 80$ mm/s (Figure 2B). Unbiased, perfect choices would result in a step function, which is approached by the psychometric curve obtained when the first stimulus $S_A = 80$ mm/s. If instead the second stimulus is equal to the mean stimulus, $S_B = 80$ mm/s, performance is poor.

To simulate an experiment, we started both recurrent and Hebbian dynamics, Equations 1 and 2, from random initial conditions. We then delivered a sequence of stimulus pairs with the same distribution of intensities and ISIs as in the experiment with human subjects (see STAR Methods for details). After $S_A^{(n)}$ —the base stimulus in trial n —has been delivered, the external input $\vec{\zeta}(t)$ is set to 0 and $S_B^{(n)}$ —the comparison stimulus in the same trial—is delivered after a delay. If at the onset of $S_B^{(n)}$, its intensity is larger than the internal representation of the first stimulus, as encoded by the RNN, the RNN makes the choice $S_B^{(n)} > S_A^{(n)}$, and vice versa. Collecting choices in this simple way, we found clear evidence of contraction bias in the accuracy of the choices of the RNN (Figure 2E), revealing that the internal representation of the RNN contracts toward the average of past stimuli. As seen for human participants, the model’s psychometric curves show good accuracy of the RNN on trials where the first stimulus intensity is near the mean stimulus intensity, while accuracy is lower on trials where the second stimulus intensity is near the distribution mean (Figure 2F). The quantitative differences in the shape of the curves between Figures 2B and 2F may stem from several factors, including the number of subjects. However, our main focus is to highlight the qualitatively similar behavior, specifically the steeper increase in response around the mid-stimulus in one condition (dark green) compared with the other (light green).

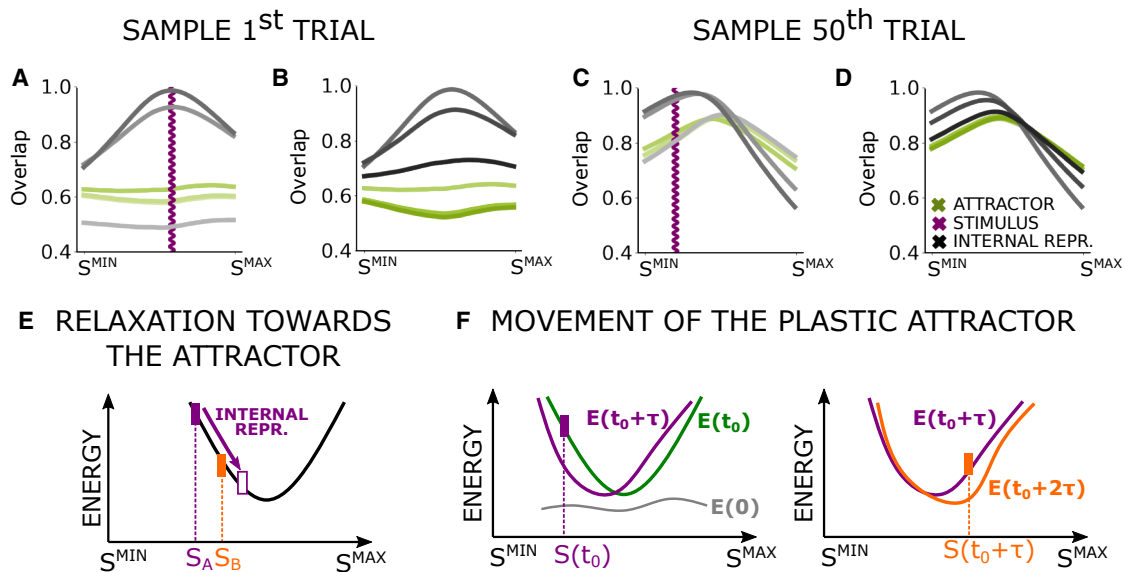


Figure 3. Dynamics of a neural network model capturing perceptual biases

Top row: RNN dynamics during the 1st and 50th trials are visualized using the overlap profile—the normalized cosine similarity between the RNN state $\vec{V}(t)$ and input patterns \vec{x}^{μ} for each stimulus intensity S^{μ} . Snapshots at selected time steps (10, 14, 20 and 20, 36, and 49 of a reference memory trial) show the evolution during and after stimulation (gray lines darken with time). Green lines depict the overlap profile of the attractor Equation 1, with weights fixed at $\vec{W}(t)$. Violet bars mark stimulus intensity. See Videos S1 and S2 for full trial dynamics.

(A) The RNN state aligns with the presented stimulus.

(B) After stimulus removal, the overlap flattens.

(C and D) By the 50th trial, Hebbian plasticity reshapes the connectivity, forming a peaked attractor. Bottom row: plasticity sculpts a dynamic attractor, illustrated via a sketch of the energy landscape near its minimum, across network states V representing different stimulus strengths.

(E) After stimulus S_A (filled violet), the RNN relaxes toward an attractor, whose intensity is the trace of past stimuli. In a working memory task, a later stimulus S_B (orange) may be misclassified if the memory of S_A (empty violet) has shifted beyond S_B , leading to contraction bias.

(F) Sketched energy profiles over time: initially (gray), no bias is present. Repeated stimulation induces a minimum that shifts progressively due to Hebbian plasticity. Effects in (E) and (F) occur concurrently.

History dependence of choices in working memory task

A fundamental feature of contraction bias reported in previous studies^{3,32} is recency—the dependence of choices on the stimuli of the immediately preceding trial. We uncover this effect in the current dataset in Figure 2C. The colored lines show the probability of making the choice $S_B^{(n)} > S_A^{(n)}$ for the three different pairs in the n th trial as a function of the second stimulus in the $(n - 1)$ th trial, $S_B^{(n-1)}$. For all three pairs of stimuli shown in Figure 2C, the weaker the preceding stimulus $S_B^{(n-1)}$, the higher the probability that $S_B^{(n)}$ was correctly classified as stronger than $S_A^{(n)}$. This shows that the stimulus $S_B^{(n-1)}$ is strongly weighted in the continuously updating average toward which $S_A^{(n)}$ is contracted. Tracing the influence of the six most recently received stimuli, i.e., $(S_B^{(k)}, S_A^{(k)})$ for $k = n, n - 1$ and $n - 2$, on the choice at trial n using the coefficient values β of a generalized linear model (GLM) reveals a positive effect of the last stimulus (meaning that subjects correctly used the current trial's S_B intensity to make choices) but a negative effect of previous stimuli; this effect decays for stimuli that are further back in time (see Figure 2D). This is replicated in the model. As shown in Figure 2G, choices made based on RNN representations have the same history dependence on previous choices as do the choices of human

subjects, which is further confirmed by a GLM analysis of RNN-based choices (Figure 2H).

In appendix working memory paradigm of Akrami et al.³, we also reproduce the experimental results from a published dataset by Akrami et al.³ which corresponds to a different working memory experiment involving a distinct sensory modality, an auditory delayed comparison task. Specifically, we reproduce the contraction bias and history dependence of choice observed in that experimental design, as well as the relationship between contraction bias and the time interval between the two stimuli.

Recurrent dynamics and continuous plasticity yield perceptual biases

How do the perceptual biases arise in the RNN model? As shown in the Videos S1 and S2 and captured in the snapshots in Figures 3A–3D, at the beginning of a session (Figure 3A), when the connectivity is still close to its random initialization, the RNN develops a large overlap with its first external input (violet line) as seen by the overlap profiles representing the cosine similarity of a given network state with all the different input patterns (going from light gray to dark gray as time increases). However, as soon as the external input is released, the RNN quickly decorrelates and the overlap profile becomes almost flat (Figure 3B). In other words, no particular stimulus intensity is represented by the network. After a few more trials, the effect of Hebbian

plasticity becomes apparent: overlap profiles remain peaked even after the stimulus offset, meaning that the network represents a uniquely identifiable stimulus intensity every time, the one corresponding to the peak of the black bell-shaped curve in Figure 3D. This behavior reflects the emergence of an attractor state with a peaked overlap profile. For *binary* morphing patterns and one-shot learning of the connectivity, Blumenfeld et al.³¹ indeed proved that a unique attractor emerges. In our case, the attractor emerges from the interplay of continuous Hebbian plasticity, the morphing input pattern structure, and intermittent stimuli presented, making it harder to analyze. However, we can make the attractor explicit at each time step t by freezing the connections $\bar{W}(t)$ and letting the network converge to the fixed point $\bar{V}^{FP}(t)$ of the recurrent dynamics (Equation 1). The overlap profiles of the attractors at the sample frozen steps are shown in green in Figures 3A–3D (light to dark as time increases); see also Videos S1 and S2.

The attractor constrains the dynamics of the network. The movement of the attractor, due to the Hebbian plasticity and the incoming stimuli, leads to choices that exhibit perceptual biases. The argument can be simplified by first considering a network with fixed weights at a certain time t . The position of the attractor associated with these weights, and thus the stimulus intensity it encodes, can be viewed as the minimum of an effective energy function $E(t)$ in the vicinity of its minimum, as sketched in Figures 3E and 3F. We do not provide an explicit form of the energy function, as is done for simpler models.^{23,33} Instead, our goal is to offer a conceptual sketch that aids understanding, where the attractor is envisioned as the bottom of a valley, with the dynamics of the RNN following the contours of this valley. At the offset of the first stimulus S_A in a working memory trial (violet in Figure 3E), the state of the RNN evolves to minimize the energy function. The stimulus intensity encoded by the RNN relaxes toward that corresponding to the attractor. In the trial sketched in Figure 3E, this relaxation reduces the distance between the internal representation of the first stimulus and the intensity of the second stimulus, S_B (shown in orange). If the ISI is long enough for the internal representation to “overtake” the intensity of the second stimulus, the RNN makes an erroneous choice by classifying the intensity of the second stimulus as lower than the first.

Although the attractor’s location will generally be close to the mean stimulus intensity, ongoing Hebbian plasticity causes its position to be continuously influenced by the external stimuli. We refer to the *trace of past stimuli* as the quantity tracked by the attractor’s location. While it acts as a sort of recency-weighted average of past stimuli, the precise relationship between past stimuli and attractor location is complex and nonlinear, as detailed in the discussion. Figure 3F illustrates the effect of two stimuli delivered one after the other, the first weaker and the second stronger than the attractor minimum. At the beginning of a session at $t = 0$, the energy $E(0)$ is almost flat. Once a sequence of trials has taken place ($t = t_0$), the energy has a definable minimum emerging through Hebbian learning. If a stimulus is then delivered (violet), this induces a change in energy, which we illustrate in the energy profile $E(t_0 + \tau)$ after some time τ . The same occurs with consequent stimuli (orange). The movement of the attractor leads to the his-

tory dependence of the choice. In the working memory experiment, if the second stimulus of the $(n - 1)$ th trial was strong, the attractor will have shifted toward higher stimulus intensities; so the internal representation of the first stimulus in the next trial will relax toward the new attractor and might overtake the second stimulus, $S_B = 80$, leading to the patterns of misclassification seen in Figure 2G. Recall that the accuracy of the correct choice for (70, 80) was high if $S_B^{(n-1)}$ was weak, and accuracy declined for progressively stronger values of $S_B^{(n-1)}$. If, instead, the second stimulus of the $(n - 1)$ th trial was weak, the attractor will have shifted toward lower stimulus intensities, facilitating the correct classification of $S_B^{(n)} > S_A^{(n)}$.

Reference memory

In this section, we show that the concurrent dynamical processes in the RNN model—the movement of the attractor and the relaxation of the RNN representation toward it—can also account for the repulsive biases observed in a reference memory task in humans and rats, without any tuning of the hyper-parameters $\gamma, \kappa, a, \lambda$. In particular, we show that the serial dependence of choices in *contraction* bias is (perhaps counter-intuitively) an instance of *repulsive* bias.

Experimental data are from the reference memory study of Hachen et al.,¹⁰ where the rat was tasked with reporting whether each single presented stimulus was weak or strong across a series of trials (Figure 1A). Since Hachen et al.¹⁰ showed that in well-trained rats, the internal representation of the threshold θ separating weak from strong stimuli was primarily influenced by the distribution of stimuli, rather than the stimulus/reward contingency, we simulated the task by driving the RNN using stimulus distributions akin to the experimental one, without modeling reward. The model’s choice occurs at the stimulus offset: if the internal representation, i.e., the peak of the overlap profile of the RNN, relaxes toward higher intensities, the delivered stimulus is read out as weak; whereas if it relaxes toward lower intensities, the delivered stimulus is read out as strong.

Hachen et al.¹⁰ showed that rats internally construct a representation of θ to categorize stimuli over successive trials and do not reliably retain it across sessions; the psychometric curves derived from the first three trials of each session deviate only marginally from chance level, as compared with those derived from three random trials (Figure 4A). This poor initial performance highlights the importance of learning the network connectivity *dynamically* for any neural network model of this behavior. Indeed, an initial set of trials is also necessary in our RNN model in order to form an attractor that is correlated with past inputs. Consequently, during the first three trials, neuronal configurations relax toward states that are not yet correlated with the effective inputs, resulting in decision probabilities closer to chance level (Figure 4E).

The repulsive bias in the reference memory task can be seen in two ways. Figure 4B shows the psychometric curves of rats that were exposed to stimulus distributions biased toward weak and strong stimuli (small light and dark blue histograms, respectively). The psychometric curves in Figure 4B show that an increased frequency of strong stimuli shifts the point of subjective equality (PSE) toward stronger stimuli, and vice versa for an increased frequency of weak stimuli. The same

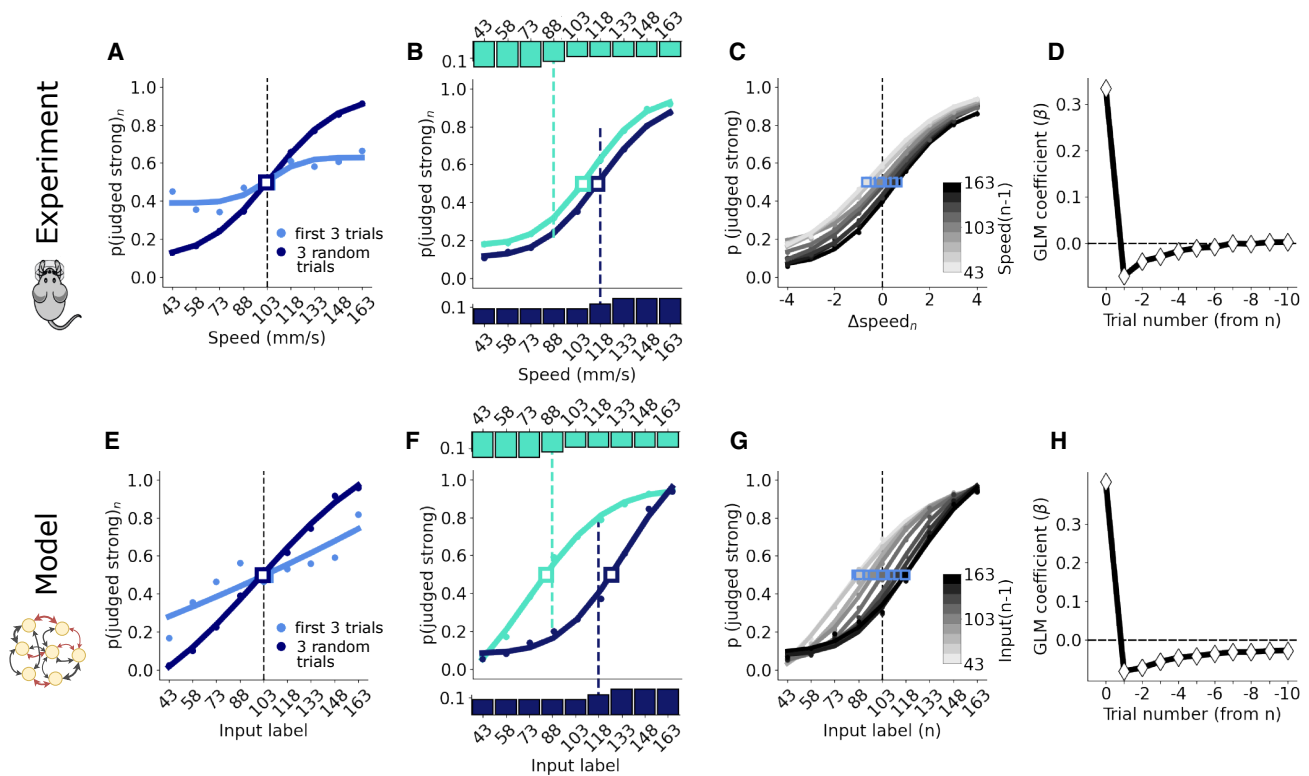


Figure 4. Psychophysical and RNN results in a reference memory task

In (A)–(D), the experimental data from Hachen et al.¹⁰ are replotted.

(A) Psychometric curves obtained from the first three trials of each session (light blue) and from three randomly selected trials per session (dark blue). The dashed vertical lines in (A)–(C) mark the category boundary. The point of subjective equality (PSE) is denoted by a square.

(B) Psychometric curves for sessions where input distributions had an average stimulus intensity of 88 (light blue-green histogram) and 118 (dark blue histogram). The input stimuli “attract” the psychometric curves toward the mean of the distribution. This curve shift does not occur if the distribution of the inputs is kept uniform, but the reward is delivered according to a category boundary displaced from the mean of the distribution (not reported here, see Hachen et al.¹⁰).

(C) The probability of categorizing a stimulus as “strong” conditioned on the intensity of the previous stimulus. Blue squares represent the PSE.

(D) Coefficient values β derived from the generalized linear model (GLM) predicting choice, given the stimulus intensities of the 10 preceding trials.

(E–H) Showing results from the RNN under conditions mimicking those of the experiments. In (A) and (B) and (E)–(G), points correspond to experimental data, solid lines to sigmoidal fits (see STAR Methods).

occurs considering the psychometric curves conditioned on the $(n - 1)$ th stimulus; Figure 4C shows that the PSE of the n th trial is pushed up to stronger stimuli after experiencing a strong stimulus on the $(n - 1)$ th trial, and vice versa. The same bias emerges in the RNN model as a consequence of Hebbian plasticity. High-intensity inputs shift the location of the attractor toward high-intensity inputs, cf. Figure 3F. This shift translates to the effect shown in both measurements: a higher probability of categorizing a given stimulus as weak (or strong), for a stimulus distribution biased toward higher (or lower) intensities, as shown in Figure 4F, and a shift in the PSE of the psychometric curves, given the different $(n - 1)$ th stimuli (Figure 4G). In appendix [asymmetric shift of psychometric curves in Hachen et al.¹⁰](#) and [individual differences](#), we discuss the quantitative disparity between Figures 4B vs. 4F, which appears to stem from individual subject differences—a phenomenon we can also replicate with different random realizations. Here, we emphasize that the results presented reveal that the model captures, in general, the relative importance of the preceding stimuli on the choices, as reported in Figures 4D and 4H.

The one-back experimental paradigm—Plasticity and relaxation at a glance

Our analysis highlighted the importance of the movement of the attractor in the RNN for recency and repulsive biases to emerge in perception. To challenge the generality of our model, we designed a novel one-back task. Participants undertake a sequence of trials, each of which presents a single tactile stimulus, similar to the reference memory task. However, triggered by the go cue, they must compare the current trial’s stimulus not to a fixed category boundary but to the stimulus of the previous trial. Comparing the current stimulus to the preceding stimulus is akin to the working memory task; however, there are marked distinctions. First, the two stimuli are not bundled together in one trial but are distributed to sequential trials. Second, each stimulus has two functions: it is initially the comparison stimulus for judgment relative to the preceding trial, but with choice made, it now becomes the base stimulus for the next trial (Figure 5A). A crucial feature of the paradigm that tests for the dynamics of the attractor is that stimulus intensities are grouped into two clouds, a “low”-intensity cloud with stimuli

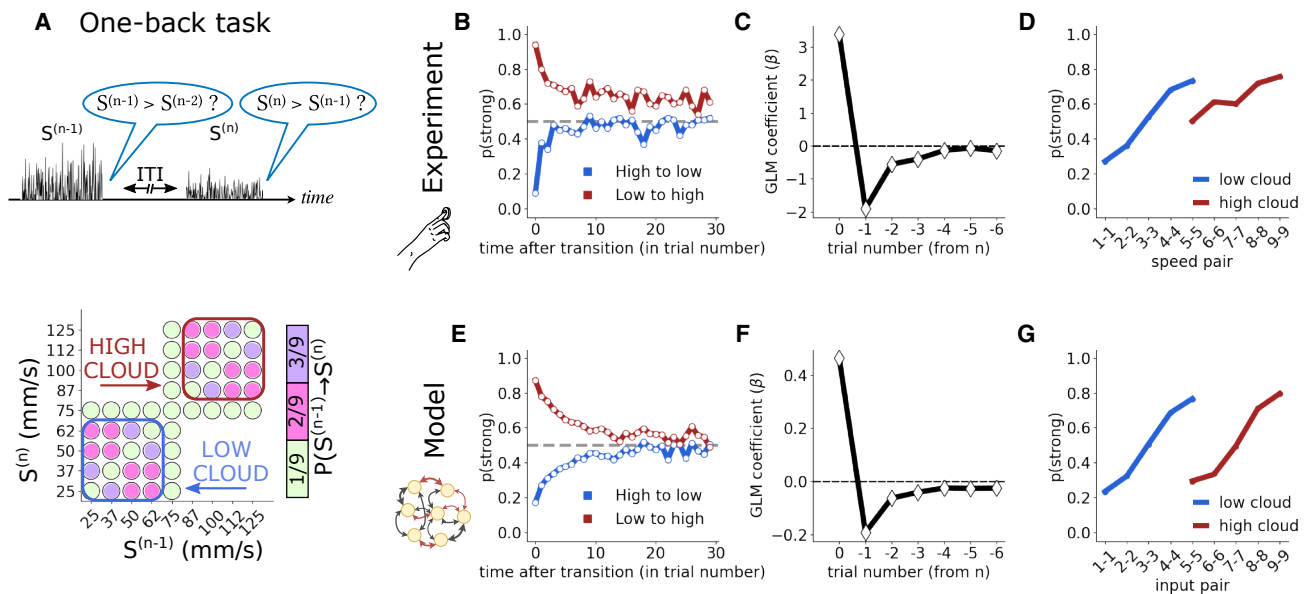


Figure 5. The one-back task highlights the plasticity of the attractor as it moves between stimulus clouds

(A) We introduce the novel one-back task, where subjects are tasked to compare the strength of one stimulus per trial with the strength of the stimulus in the preceding trial (top). We design the transition probabilities governing the strength of the n th stimulus given the strength of the $(n - 1)$ th stimulus for this task (bottom) such that long consecutive sequences of stimuli from the low-intensity or high-intensity “clouds” appear.

(B) Probability of human subjects categorizing a stimulus as stronger than the preceding stimulus within each cloud as a function of the time elapsed since the transition between clouds.

(C) Coefficient values (β) of the generalized linear model (GLM) of the choice at trial n , given the six preceding stimuli.

(D) Probability of categorizing $S^{(n)} > S^{(n-1)}$ when $S^{(n)} = S^{(n-1)}$. If the pair (75, 75) is preceded and followed by stimuli from the low cloud, it is considered part of the low cloud, and vice versa for the high cloud.

(E–G) Results obtained from the RNN model.

labeled from 1 – 4 and a “high” cloud of stimuli labeled from 6 – 9. After a stimulus is drawn from the low/high cloud, the next stimulus is drawn either from the same cloud or is the intermediate stimulus 5. A transition from one to the other cloud may occur only with a trial at the intermediate stimulus. The transition probabilities $p(S^{(n)}|S^{(n-1)})$, shown in Figure 5A, can lead to long sequences of stimuli constrained to one of the clouds (mean within-cloud string length is 18 trials; median is 13). This paradigm tests the prediction that subjects’ psychophysical results will reflect the attractor’s continuous movement between the two clouds, as the trace of past stimuli transitions between low and high.

The experiment comprised twenty-four human participants, each completing 750 trials. Figure 5B shows the probability that a given stimulus was judged as stronger than the previous as a function of time after a jump between clouds. By objective physical measures, within each cloud the (n) th stimulus is equally likely to be weaker or stronger than the $(n - 1)$ th stimulus. But perceptual biases were uncovered in actual participants. After a switch from the high- to the low-intensity cloud, participants consistently underestimated the strength of the (n) th stimulus (against the $(n - 1)$ th stimulus) over the first few trials (blue line), before stabilizing around the correct value of 0.5. Conversely, after switching from low- to high-intensity stimuli, stimulus intensities were overestimated initially (red). The RNN model generates analogous behavior (Figure 5E). The

model’s results are explained by the attractor migrating gradually, after the cloud switch, from its previous position around the average intensity of the opposite cloud. As the $(n - 1)$ th stimulus representation is pulled toward the earlier cloud, a strong bias produces the observed misclassifications. As the stimulus string remains in the current cloud, the attractor slowly settles to a new globally stable position with small dynamical movements around the average intensity value of the cloud.

Averaging over time, the overall influence of previous stimuli on the choice can be estimated in terms of GLM coefficients as done in the previous two paradigms, showing similar results; the choice at trial n depends, with a decaying weight, on the previous stimuli as found in both experimental and model data in Figures 5C and 5F).

The one-back paradigm highlights that the shifting of internal representations, which is responsible for contraction bias, does not occur toward the global average of stimulus intensities, which would be the center of the stimulus range, intensity 5; instead, internal representations contract toward an intensity value that changes in time. The contraction bias in the one-back task is most evident considering choices of participants after experiencing the same input successively, $S^{(n)} = S^{(n-1)}$. The probability of categorizing $S^{(n)} > S^{(n-1)}$ deviates from 0.5 in a discontinuous trend, shown in Figure 5D for human subjects. For instance, within a stream of low-intensity stimuli, the second presentation of stimulus of intensity 5 is consistently categorized

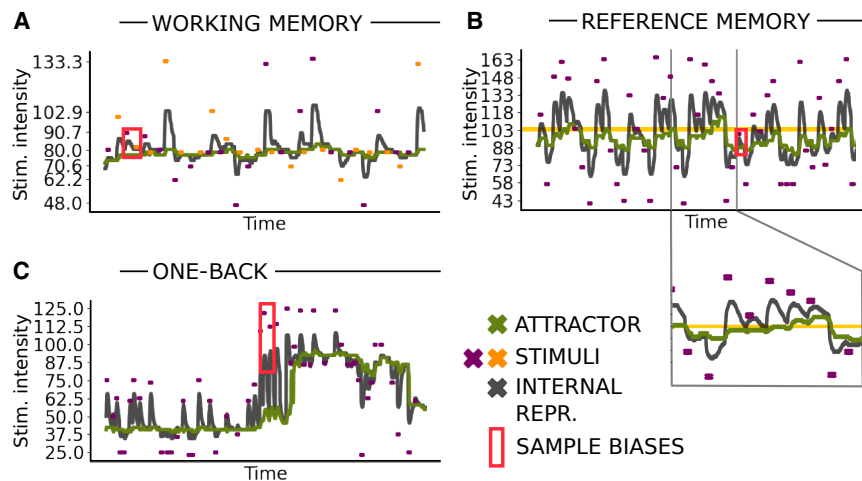


Figure 6. A single recurrent neural network with continuous plasticity reproduces two different biases across three paradigms

The dynamics driving perceptual biases across the three different experimental paradigms. Evolution of the network in a sample working memory (A), reference memory (B), one-back (C) task.

The external stimulus is depicted in violet and in orange when corresponding to S_B in (B). The internal representation (dark gray) and attractor position (green) at each time step is read out from the maximum of the overlap profile. The red rectangles indicate instances of incorrect choices, while the orange horizontal line in (B) represents the overall average.

as stronger than the previous stimulus of intensity 5, while the opposite is true within a stream of high-intensity stimuli. The internal representation of intensity 5 contracts toward a lower value when it is part of a low-intensity sequence, and vice versa, consistently biasing the categorization of the second instance of intensity 5. The RNN model fully reproduces this behavior (Figure 5G). Indeed, during the ITI, the internal representation of the stimulus relaxes toward the running attractor, replicating the temporal dynamics of actual human subjects. Note that the asymmetry in the experimental results for stimuli from the high vs. the low cloud in Figure 5D and the decay of the red curve to a value larger than 0.5 in Figure 5B arise from a combination of the Weber-Fechner law and the stochasticity in the stimulation setup. In our model, where the Weber-Fechner effect is not included, the asymmetries in Figures 5E and 5G stem solely from this stochasticity in the stimulation, as detailed in appendix [the effect of the stochasticity in the stimulation setup](#). We emphasize that, again, this simulation was achieved with no external re-setting of the hyper-parameters γ, λ, η, g used previously for working memory and reference memory simulations.

DISCUSSION

Perceptual biases^{1,3,5,7–11} offer a window into the computational processes underlying perceptual judgments, but the neural mechanisms from which they emerge remained poorly understood. Here, we proposed a simple RNN model with ongoing Hebbian plasticity that reproduces perceptual biases across four datasets spanning three different experimental paradigms, without any per-paradigm fine-tuning. Combining Hebbian plasticity with a set of biologically inspired input patterns yields a dynamic point attractor. Contraction arises from relaxation toward this plastic attractor, and the movement of the attractor, driven by Hebbian plasticity, leads to both repulsion and recency effects. Our newly designed one-back task accentuates the movement of this plastic attractor, and the RNN model correctly predicted the novel bimodal contraction bias we observed in the performance of human subjects on the task.

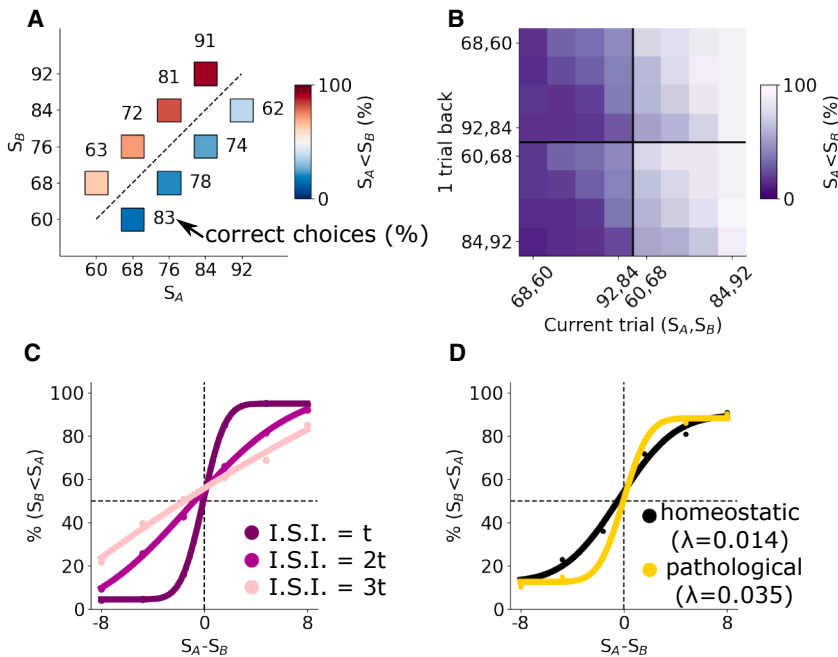
Figure 6 provides a closer look at the time traces of the neurons in our model. We report the intensity values of the internal repre-

sentation (gray) and of the attractor (green) in sample simulations per paradigm. While the dynamics of neurons and weights follow the same dynamical Equations 1 and 2 in all three settings, it is the task-specific input distribution that leads to different choice statistics. Wrong choices, highlighted by red rectangles, reproduce task-specific perceptual biases.

There are two important timescales in our model. One governs the relaxation of the internal representation in the absence of stimulation, while another governs the movement of the plastic attractor. While the relaxation occurs on a shorter timescale than the movement of the attractor, the two timescales are comparable, and both evolve in time. The movement of the attractor's location, which we refer to as the trace of past stimuli, depends on the states of the RNN $\vec{V}(t)$, which in turn depend on the attractor itself and on the external stimuli. These timescales are thus only indirectly fixed by the hyper-parameters of the model; instead, the memory decay of a stimulus representation occurs over an *effective timescale* that emerges from the interaction of individual neurons and from the continuously evolving connectivity. We show in Figure S5 that this timescale is indeed much longer than the timescale of individual neurons γ , meaning our model addresses the key challenge of modeling working memory, namely, overcoming the discrepancy between the timescale of the dynamics of individual neurons and the timescale of memory.¹⁷ This is particularly evident in the one-back paradigm sample reported in Figure 6C; a given stimulus influences the speed of movement and the effective value of both the internal representation and the attractor differently, depending on whether it occurs at the beginning or end of a period of continuous stimulation with stimuli from one of the two clouds. A promising direction for future research would be to explore how our findings relate to previous studies suggesting that the perception of probabilities extends beyond the running average.^{34,35}

The idea that perceptual biases may arise from attractive dynamics has long been hypothesized. Initially, this concept was framed through a phenomenological model involving two interacting attractors with different timescales.^{10,14,15} More recently, Boboeva et al.¹⁹ incorporated this idea into a computational model that postulates two static continuous attractor neural networks with predetermined connectivity and tuning curves to

MODEL RESULTS ON THE EXPERIMENTAL PARADIGM of (Akrami, et. al., 2018, Nature)



reproduce a single experiment. In our model, we let a single RNN learn its connectivity from stimuli via a local, biologically inspired learning rule, and we obtained a model that generalizes across behavioral tasks and species. While in Boboeva et al.¹⁹ contraction bias results from the difference in timescale of the two networks and the presence of adaptation only in one of them, in our model, contraction bias is an emergent effect of learning.

Our model also successfully reproduces a key result that was previously considered evidence for the concept that perceptual biases arise from the interaction between two neural networks; in a study by Akrami et al.,³ inactivation of the posterior parietal cortex (PPC) subregion was shown to reduce contraction bias in an auditory working memory tasks, despite no loss of memory capabilities. We found that our RNN model replicates this behavior following PPC inactivation by adjusting synaptic homeostasis, as shown in Figure 7. Modifying the synaptic homeostasis effectively places the model in a regime in which the attractor tends to remain stuck at the intensity of the last perceived stimulus. This adjustment preserves memory capacity while reducing contraction bias, as the relaxation dynamics only minimally varies the effective coded intensity. Our findings suggest an alternative role for the PPC subnetwork in the modulation of contraction bias, see appendix “Increasing synaptic homeostasis reproduces the effect of optogenetic inactivation of PPC”.

In the future, it will be intriguing to confront our high-dimensional RNN model with high-density recordings from brain regions with recurrent circuitry. First, it would be interesting to look for signatures of a dynamically moving attractor in neural recordings from animals performing perceptual tasks, for example,

Figure 7. Increasing synaptic homeostasis reproduces the effect of optogenetic inactivation of PPC. Replication of the experimental results of Akrami et al.³ with our RNN model see appendix working memory paradigm of Akrami et al.³

(A–C) reproduce Figure 1B, Figures 2A–2C, Figure 1D, and Figure 3C of Akrami et al.³ (A) Probability of correct choices for eight stimulus pairs and (B) depending on previous pair. (C) Performance across a subset of pairs depending on the ISI intervals.

(D) Performance variations when $\lambda = 0.035$ (yellow). Solid lines in (C) and (D) are fits.

by analyzing how stimulus intensities are represented by the neural population and how neural activity relaxes after stimulus removal. More generally, it is interesting to contrast neural representations and their geometry^{36–38} found *in vivo* with representations obtained from RNN models with continuous Hebbian plasticity and RNNs trained with gradient descent.

Our model has some limitations. We assumed a dense network with all-to-all recurrent connections, although it is unclear whether this accurately reflects

the connectivity in the brain’s memory storage networks, such as vibrissal motor cortex (vM1).³⁰ Additionally, our model assumes that inputs predominantly originate from a single region, with a balanced input distribution (50% code external input intensity positively, and 50% code it negatively). However, recordings from VS1 suggest that most units actually code intensity positively. Future work should investigate the effects of sparser connectivity and unbalanced input schemes on the network’s behavior.

Our model shows how attractors, which have long been contemplated to code intensity in the brain,^{10,18,19} can be learnt continuously from the input distribution using a local learning rule. Our study also highlights the potential significance of transient RNN dynamics for encoding and processing information. The mathematical analysis of an RNN with ongoing Hebbian plasticity driven by intermittent, correlated inputs is a significant challenge. The recent dynamical mean-field theory for coupled neuronal-synaptic dynamics by Clark and Abbott³⁹ represents an excellent starting point. Finally, extending this model beyond sensory perception is another promising direction for further research.

RESOURCE AVAILABILITY

Lead contact

Further information and requests for resources and reagents should be directed to and will be fulfilled by the lead contact, Sebastian Goldt (sgoldt@sissa.it).

Materials availability

This study did not generate new, unique reagents.

Data and code availability

- Data reported in this paper are available from the [lead contact](#) upon reasonable request.
- The code used to simulate the RNN is available at <https://osf.io/msbuh/>, doi:10.17605/OSF.IO/MSBUH.
- Any additional information required to reanalyze the data reported in this paper is available from the [lead contact](#) upon request.

ACKNOWLEDGMENTS

We thank Iacopo Hachen for valuable discussions and for providing the data used in [Figure 4](#). We thank Alessandro Treves for valuable discussions on Blumenfeld et al.³¹ F.S. is supported by the QBio Junior Research Chair program of the QBio initiative of ENS-PSL. M.D. acknowledges financial support from the Human Frontier Science Program, project RGP0017/2021 and Italian Ministry PRIN 2022 contract 20224FWF2. S.G. gratefully acknowledges funding from the European Research Council (ERC) for the project “beyond2” under the European Union’s Horizon 2020 research and innovation programme, grant agreement ID 101166056, and co-funding from Next Generation EU in the context of the National Recovery and Resilience Plan, Investment PE1 – Project FAIR “Future Artificial Intelligence Research.” This resource was co-financed by the Next Generation EU [DM 1555 del 11.10.22].

AUTHOR CONTRIBUTIONS

F.S., M.D., and S.G. designed the study. F.S. proposed the theoretical model, developed the simulations, and performed the experimental-theoretical comparisons. D.G. performed the experiments on the working memory paradigm, and Y.C. performed the experiments on the one-back paradigm. F.S. and S.G. wrote the first draft of the paper, and F.S., S.G., and M.D. carried out subsequent revisions. All authors approved the final manuscript.

DECLARATION OF INTERESTS

The authors declare no competing interests.

STAR★METHODS

Detailed methods are provided in the online version of this paper and include the following:

- [KEY RESOURCES TABLE](#)
- [EXPERIMENTAL MODEL AND STUDY PARTICIPANT DETAILS](#)
 - Human working memory task
 - Human one-back task
- [METHOD DETAILS](#)
 - Experimental protocols
 - RNN model
 - Synaptic Weight Changes and Biological Plausibility
 - Increasing synaptic homeostasis reproduces the effect of optogenetic inactivation of PPC³
 - Asymmetric shift of psychometric curves in Hachen et al.¹⁰ and individual differences
 - The effect of the stochasticity in the stimulation setup
 - Further details on time scales in the model

SUPPLEMENTAL INFORMATION

Supplemental information can be found online at <https://doi.org/10.1016/j.neuron.2025.09.037>.

Received: August 13, 2024

Revised: December 23, 2024

Accepted: September 25, 2025

Published: November 5, 2025

REFERENCES

1. Ashourian, P., and Loewenstein, Y. (2011a). Bayesian inference underlies the contraction bias in delayed comparison tasks. *PLoS One* 6, e19551. <https://doi.org/10.1371/journal.pone.0019551>.
2. Fassihi, A., Akrami, A., Esmaili, V., and Diamond, M.E. (2014). Tactile perception and working memory in rats and humans. *Proc. Natl. Acad. Sci. USA* 111, 2331–2336. <https://doi.org/10.1073/pnas.1315171111>.
3. Akrami, A., Kopec, C.D., Diamond, M.E., and Brody, C.D. (2018). Posterior parietal cortex represents sensory history and mediates its effects on behaviour. *Nature* 554, 368–372. <https://doi.org/10.1038/nature25510>.
4. Tal-Perry, N., and Yuval-Greenberg, S. (2022). Contraction bias in temporal estimation. *Cognition* 229, 105234. <https://doi.org/10.1016/j.cognition.2022.105234>.
5. Hollingworth, H.L. (1910). The central tendency of judgment. *J. Philos. Psychol. Sci. Methods* 7, 461–469. <https://doi.org/10.2307/2012819>.
6. Joynson, R.B., Newson, L.J., and May, D.S. (1965). The limits of over-constancy. *Q. J. Exp. Psychol.* 17, 209–216. <https://doi.org/10.1080/17470216508416434>.
7. Gibson, J.J., and Radner, M. (1937). Adaptation, after-effect and contrast in the perception of tilted lines. I. Quantitative studies. *J. Exp. Psychol.* 20, 453–467. <https://doi.org/10.1037/h0059826>.
8. Magnussen, S., and Johnsen, T. (1986). Temporal aspects of spatial adaptation. A study of the tilt aftereffect. *Vision Res.* 26, 661–672. [https://doi.org/10.1016/0042-6989\(86\)90014-3](https://doi.org/10.1016/0042-6989(86)90014-3).
9. Fritsche, M., Mostert, P., and de Lange, F.P. (2017). Opposite effects of recent history on perception and decision. *Curr. Biol.* 27, 590–595. <https://doi.org/10.1016/j.cub.2017.01.006>.
10. Hachen, I., Reinartz, S., Brasselet, R., Stroligo, A., and Diamond, M.E. (2021). Dynamics of history-dependent perceptual judgment. *Nat. Commun.* 12, 6036. <https://doi.org/10.1038/s41467-021-26104-2>.
11. Karim, M., Harris, J.A., Langdon, A., and Breakspear, M. (2013). The influence of prior experience and expected timing on vibrotactile discrimination. *Front. Neurosci.* 7, 255. <https://doi.org/10.3389/fnins.2013.00255>.
12. Fritsche, M., Spaak, E., and De Lange, F.P. (2020). A Bayesian and efficient observer model explains concurrent attractive and repulsive history biases in visual perception. *eLife* 9, e55389. <https://doi.org/10.7554/eLife.55389>.
13. Oña-Jodar, T., Prat-Ortega, G., Li, C., Dalmau, J., Compte, A., and Rocha, J.D.L. (2023). NEURAL MECHANISMS UNDERLYING WORKING MEMORY ERRORS. *IBRO Neurosci. Rep.* 15, S672–S673. <https://doi.org/10.1016/j.ibneur.2023.08.1351>.
14. Yousefi Darani, Z.Y., Hachen, I., and Diamond, M.E. (2023). Dynamics of the judgment of tactile stimulus intensity. *Neuromorph. Comput. Eng.* 3, 014014. <https://doi.org/10.1088/2634-4386/acc08e>.
15. Diamond, M.E., and Toso, A. (2023). Tactile cognition in rodents. *Neurosci. Biobehav. Rev.* 149, 105161. <https://doi.org/10.1016/j.neubiorev.2023.105161>.
16. Hahn, M., and Wei, X.-X. (2024). A unifying theory explains seemingly contradictory biases in perceptual estimation. *Nat. Neurosci.* 27, 793–804. <https://doi.org/10.1038/s41593-024-01574-x>.
17. Seung, H.S. (1996). How the brain keeps the eyes still. *Proc. Natl. Acad. Sci. USA* 93, 13339–13344. <https://doi.org/10.1073/pnas.93.23.13339>.
18. Machens, C.K., Romo, R., and Brody, C.D. (2005). Flexible control of mutual inhibition: a neural model of two-interval discrimination. *Science* 307, 1121–1124. <https://doi.org/10.1126/science.1104171>.
19. Boboeva, V., Pezzotta, A., Clopath, C., and Akrami, A. (2023). From recency to central tendency biases in working memory: a unifying network model. *eLife* 12, RP86725. <https://doi.org/10.7554/eLife.86725.2>.
20. Grossberg, S. (1969). On learning, information, lateral inhibition, and transmitters. *Math. Biosci.* 4, 255–310. [https://doi.org/10.1016/0025-5564\(69\)90015-7](https://doi.org/10.1016/0025-5564(69)90015-7).

21. Hopfield, J.J. (1984). Neurons with graded response have collective computational properties like those of two-state neurons. *Proc. Natl. Acad. Sci. USA* *81*, 3088–3092. <https://doi.org/10.1073/pnas.81.10.3088>.
22. Sompolinsky, H., Crisanti, A., and Sommers, H.J. (1988). Chaos in random neural networks. *Phys. Rev. Lett.* *61*, 259–262. <https://doi.org/10.1103/PhysRevLett.61.259>.
23. Treves, A. (1990). Graded-response neurons and information encodings in autoassociative memories. *Phys. Rev. A* *42*, 2418–2430. <https://doi.org/10.1103/physreva.42.2418>.
24. Schönsberg, F., Monasson, R., and Treves, A. (2024). Continuous quasi-attractors dissolve with too much- or too little - variability. *PNAS Nexus* *3*, pgae525. <https://doi.org/10.1093/pnasnexus/pgae525>.
25. Hebb, D. (1949). *The Organization of Behavior: A Neuropsychological Theory* (John Whisking Wiley).
26. Oja, E. (1982). A Simplified neuron model as a principal component analyzer. *J. Math. Biol.* *15*, 267–273. <https://doi.org/10.1007/BF00275687>.
27. Sussillo, D., and Abbott, L.F. (2009). Generating coherent patterns of activity from chaotic neural networks. *Neuron* *63*, 544–557. <https://doi.org/10.1016/j.neuron.2009.07.018>.
28. Wallace, M.L., and Sabatini, B.L. (2023). Synaptic and circuit functions of multitransmitter neurons in the mammalian brain. *Neuron* *111*, 2969–2983. <https://doi.org/10.1016/j.neuron.2023.06.003>.
29. Arabzadeh, E., Petersen, R.S., and Diamond, M.E. (2003). Encoding of whisker vibration by rat barrel cortex neurons: implications for texture discrimination. *J. Neurosci.* *23*, 9146–9154. <https://doi.org/10.1523/JNEUROSCI.23-27-09146.2003>.
30. Fassihi, A., Akrami, A., Pulecchi, F., Schönfelder, V., and Diamond, M.E. (2017). Transformation of perception from sensory to motor cortex. *Curr. Biol.* *27*, 1585–1596. <https://doi.org/10.1016/j.cub.2017.05.011>.
31. Blumenfeld, B., Preminger, S., Sagi, D., and Tsodyks, M. (2006). Dynamics of memory representations in networks with novelty-facilitated synaptic plasticity. *Neuron* *52*, 383–394. <https://doi.org/10.1016/j.neuron.2006.08.016>.
32. Raviv, O., Ahissar, M., and Loewenstein, Y. (2012). How recent history affects perception: the normative approach and its heuristic approximation. *PLoS Comput. Biol.* *8*, e1002731. <https://doi.org/10.1371/journal.pcbi.1002731>.
33. Amit, D.J., Gutfreund, H., and Sompolinsky, H. (1985). Storing infinite numbers of patterns in a spin-glass model of neural networks. *Phys. Rev. Lett.* *55*, 1530–1533. <https://doi.org/10.1103/PhysRevLett.55.1530>.
34. Babayan, B.M., Uchida, N., and Gershman, S.J. (2018). Belief state representation in the dopamine system. *Nat. Commun.* *9*, 1891. <https://doi.org/10.1038/s41467-018-04397-0>.
35. Gallistel, C.R., Krishan, M., Liu, Y., Miller, R., and Latham, P.E. (2014). The perception of probability. *Psychol. Rev.* *121*, 96–123. <https://doi.org/10.1037/a0035232>.
36. Bernardi, S., Benna, M.K., Rigotti, M., Munuera, J., Fusi, S., and Salzman, C.D. (2020). The geometry of abstraction in the hippocampus and prefrontal cortex. *Cell* *183*, 954–967. <https://doi.org/10.1016/j.cell.2020.09.031>.
37. Chung, S., and Abbott, L.F. (2021). Neural population geometry: An approach for understanding biological and artificial neural networks. *Curr. Opin. Neurobiol.* *70*, 137–144. <https://doi.org/10.1016/j.conb.2021.10.010>.
38. Nogueira, R., Rodgers, C.C., Bruno, R.M., and Fusi, S. (2023). The geometry of cortical representations of touch in rodents. *Nat. Neurosci.* *26*, 239–250. <https://doi.org/10.1038/s41593-022-01237-9>.
39. Clark, D.G., and Abbott, L. (2024). Theory of coupled neuronal-synaptic dynamics. *Phys. Rev. X* *14*, 021001. <https://doi.org/10.1103/PhysRevX.14.021001>.
40. Schwarz, C., and Chakrabarti, S. (2016). Whisking control by motor cortex. In *Scholarpedia of Touch*, T. Prescott, E. Ahissar, and E. Izhikevich, eds. (Atlantis Press), pp. 751–769. https://doi.org/10.2991/978-94-6239-133-8_55.
41. Lee, H., and Lee, S.-H. (2023). Boundary updating as a source of history effect on decision uncertainty. *Iscience* *26*, 108314. <https://doi.org/10.1016/j.isci.2023.108314>.
42. Docherty, M., Bradford, H.F., and Wu, J.Y. (1987). Co-release of glutamate and aspartate from cholinergic and GABAergic synaptosomes. *Nature* *330*, 64–66. <https://doi.org/10.1038/330064a0>.
43. Jonas, P., Bischofberger, J., and Sandkühler, J. (1998). Corelease of two fast neurotransmitters at a central synapse. *Science* *281*, 419–424. <https://doi.org/10.1126/science.281.5375.419>.
44. Morales, M., and Margolis, E.B. (2017). Ventral tegmental area: cellular heterogeneity, connectivity and behaviour. *Nat. Rev. Neurosci.* *18*, 73–85. <https://doi.org/10.1038/nrn.2016.165>.
45. Fechner, G.T. (1860). In *Elemente der Psychophysik* (Breitkopf u. Härtel).
46. Toso, A., Fassihi, A., Paz, L., Pulecchi, F., and Diamond, M.E. (2021). A sensory integration account for time perception. *PLoS Comp. Biol.* *17*, e1008668. <https://doi.org/10.1371/journal.pcbi.1008668>.

STAR★METHODS

KEY RESOURCES TABLE

REAGENT or RESOURCE	SOURCE	IDENTIFIER
Software and algorithms		
We developed python code to simulate and analyse the recurrent neural network model.	http://osf.io/msbuh	10.17605/OSF.IO/MSBUH

EXPERIMENTAL MODEL AND STUDY PARTICIPANT DETAILS

Human working memory task

Sixteen subjects participated in the WM task. All subjects performed one WM session, lasting approximately one and a half hour for a total of 1000 trials. All subjects were volunteers and were paid after the participation, on the basis of how well they performed. The study protocol conformed to international norms and was approved by the Ethics Committee of the International School for Advanced Studies (SISSA).

Human one-back task

Twenty-four human subjects (ages 19-38), recruited from the online research participation system (SONA), participated in the One Back Experiment. The study adhered to frame of rules specified by international norms for human behaviour experiments and was approved by the Ethics Committee of the International School for Advanced Studies (SISSA).

METHOD DETAILS

Experimental protocols

Human working memory task

All experiments took place in the SENSEx Lab at SISSA. The subjects sat in front of a PC screen with their right arm resting on a pillow upon the desk, using their right index fingertip to trigger an infrared sensor in order to begin a trial; when doing so, their fingertip would be in contact with a plastic probe. The probe was attached to a shaker motor (type 4808; Bruel and Kjaer), producing vibrations in the horizontal direction, perpendicular to the fingertip. Subjects' responses were reported by pressing one of two buttons with the right hand. The experiment was automatised via LabVIEW (National Instruments).

Each stimulus was a noisy vibration, obtained by stringing together randomly sampled velocities. Velocities were obtained by sampling a normal distribution with 0 mean and defined by the standard deviation σ , ranging from 3 to 20; the mean stimulus speed is directly proportional to σ , and the motor amplifier gain was set so that the average stimulus speed would be 10 times the value of σ , measured in mm/s. The sequence of velocities formed one seed. There were 50 different seeds for each vibration mean speed, making the signature of a given vibration hard to recognise and making the mean speed a more salient feature.

The subject's finger had to remain in the sensor for the entire trial, comprising: pre-stimulus delay (500ms), base stimulus (334ms), inter-stimulus-interval (either 2, 4 or 8s), comparison stimulus (334ms), post-stimulus delay (500ms); after this sequence elapsed, a visual go cue signalled to the subject to report their choice. If the subject lost contact with the sensor before the go cue, the trial aborted, and a new one had to be initiated.

Nominal pair difficulty does not depend on the value of the stimuli, but on the relative difference between them: the closer they are, the more difficult the discrimination. We expressed this difficulty using an index, the Normalised Speed Difference (NSD), according to Weber's Law:

$$\frac{S_B - S_A}{S_B + S_A}$$

The stimuli pairs were arranged in two sets: a vertical set, in which the S_A σ was fixed at 8, and a horizontal set, in which the S_B stimulus σ was fixed at 8. For both sets, NSD values ranged from -0.25 to 0.25, including 0, and the non-fixed σ was calculated accordingly. Stimuli pairs from both sets were randomly presented during the WM sessions.

The GLM model employed a logistic link function and the parameters were estimated via Maximum Likelihood Estimation by bootstrapping trials on a subject by subject basis.

Human one-back task. The experiment was conducted in Human labs at SENSEx Lab in SISSA. The subjects were seated in a dimly lit room. The experiment was executed in LabView software. Subjects viewed a monitor screen for cues and feedback related to the task. They wore headphones that presented white noise and eliminated any external sounds from the environment or motor.

The vibrotactile stimuli were delivered to the tip of their right-hand index finger through a shaker motor (type 4808; Bruel and Kjaer). The responses were recorded using keyboard button press using left hand. The subjects saw a blue dot as fixation on the screen. They received feedback (green/red dot for correct /incorrect) on each trial through the screen.

The subjects performed a one-back memory task where the goal was to compare stimulus $S^{(n)}$ with $S^{(n-1)}$. The stimuli were generated by sampling a sequence of velocity values. The probability of velocity values was drawn from a normal distribution ($\mu = 0$, $SD = \sigma$ in mm/s). The gain of the motor was amplified 10 times. The velocity sequence was taken randomly from 50 seed values for a given trial. These were transmitted as voltage values to the shaker motor producing motion in horizontal direction using a plastic knob delivered perpendicular to the finger. The stimuli had a fixed duration of 350 ms. The judgement is based on the mean speed of the stimulus, which is subjectively reported by subjects to be “strength”, “amplitude”, or “intensity”. The mean speed of the vibration was proportional to the SD values used. Nine SD values [2.5, 3.75, 5, 6.25, 7.5, 8.75, 10, 11.25 and 12.5] were used.

Each trial began with stimulus delivery followed by response time window (timed on the screen by a graphical motion of blue colour filling a vertical bar) followed by Inter-Trial Interval (ITI). The subjects had two seconds available to report their decision by pressing keys [A or D]. After every response, the feedback was displayed on the screen. Failure to respond within the given response time resulted in a timeout. Three ITI values [0, 2 and 6 sec] were used, resulting in total ITI of 2, 4 or 8 seconds between two consecutive stimuli. Reward was delivered for correct choices (if the same stimulus intensity is presented on two consecutive trials, the correctness of the comparison is rewarded randomly.)

The experiment lasted one hour and was divided into 3 blocks (250 trials each). Each subject performed a total of 750 trials. The average performance of subjects in One Back task is 76.63 percent. All data analyses were performed using custom scripts in MATLAB. The GLM model was calculated as introduced in the previous section.

RNN model

Motivation

The four psychophysical experiments that we reproduce with our model – a reference memory, two working memory, and a one-back memory tasks – are comprised of judgements made on vibratory stimuli delivered to the fingertip of human participants, to the whiskers of the rats or by auditory stimuli delivered to rats. As far as rats and tactile stimuli are considered, the model may be considered a representation of the frontal cortical region known as vibrissal motor cortex (vM1), the main target of somatosensory cortex in rats.⁴⁰ vM1 is involved in memory, decision making, and motor planning and is a recurrent network that receives afferent inputs from the primary vibrissal somatosensory cortex (VS1).⁴⁰ In humans, the corresponding cortical regions that fulfil these tasks are currently unknown, though new evidence is emerging.⁴¹

Sources of Randomness

Different realizations of the model arise from several sources of randomness: a) The realization of the input pattern scheme; b) The realization of the stimulation protocol; c) The stochasticity in the stimulation setup for each delivered input; d) The random initialization of the weights. The random protocols for a, b, and c are described in appendix [input patterns](#), the one for d in the following [simulation](#) subsection.

Input patterns

Our model for the afferent inputs to the RNN units is inspired by physiological recordings from Arabzadeh et al.²⁹ and Fassihi et al.³⁰ of the primary vibrissal somatosensory cortex (VS1). They found that the majority of vS1 neurons in behaving rats show monotonically increasing firing rates as vibration intensity increases, with heterogeneity across neurons in regard to sensitivity (minimum vibration to cause excitation) and slope. We set the ratio of positive versus negative slopes to 50% to account for the possible presence of interneurons. Half of the neurons receive inputs whose magnitude increases with stimulus strength (positive slope), while the other half receives inputs that decreases with stimulus strength. We discretize the intensity line, from a minimum to a maximal value, into P bins labelled $\mu = 1..P$. We thus obtain an input pattern matrix $\vec{\xi}$ of size $N \times P$, where ξ_i^μ is the input unit i receives when the intensity value corresponding to the μ label is delivered to the network. We provide additional details about the generation of *input patterns* used in our model in appendix [input patterns](#). In addition, we also model the stochasticity of the experimental stimulation setup, a minor detail also explained in appendix [input patterns](#), for which, essentially, each stimulus strength is actually defined by a standard deviation. The specific values for the standard deviations used to sample vibrational stimulus strengths and their respective extraction probabilities used in both the tactile experimental sessions and the model are as follows:

Working Memory Task Standard deviation values (σ) were 4.8, 6.22222, 7.05888, 8, 9.06667, 10.28577, and 13.33333. The S_A and S_B couples were organised as depicted in [Figures 2A and 2E](#), with each pair extracted with equal probability.

Reference Memory Task Standard deviation values (σ) were 4.3, 5.8, 7.3, 8.8, 10.3, 11.8, 13.3, and 14.8. Each value was extracted with equal probability, except in [Figures 4B and 4F](#), where biased extraction probabilities were indicated in the plot.

One-Back Task Standard deviation values (σ) were 2.5, 3.75, 5, 6.25, 7.5, 8.75, 10, 11.25, and 12.5. Extraction probabilities were as reported in [Figure 5](#).

Details on the input generation. Here, we provide additional details about the generation of input patterns utilised in our model. The inputs are represented by a matrix $\vec{\xi}$ of size $N \times P$, where N is the size of the network and P are the number of bins discretising the intensity values from a minimal to a maximal one. Each row $\vec{\xi}_i$ denotes the input profile received by neuron i for every intensity value. Each column $\vec{\xi}^\mu$ represents the pattern of activity delivered to the network when a specific input intensity, labelled as μ , is present,

with ξ_i^μ denoting the intensity input received by neuron i for input μ . To discretise intensity values, we divide a range from a minimum s^{MIN} to a maximum s^{MAX} intensity into P bins, each coding for the intensity $((S^{MAX} - S^{MIN})/P)$. To generate the input patterns, we employ the following procedure: we begin with a vector of length N containing random numbers ranging from -1 to $+1$; we create a second vector by replacing each negative entry in the first vector with a random number between 0 and 1, and each positive entry with a random number between -1 and 0; we linearly interpolate $P - 2$ patterns from the first to the last vector, ensuring continuity between intensity levels; finally, we set any negative values to zero to maintain non-negativity in the generated patterns. Through this process, we obtain a set of input patterns ranging from ξ_i^0 to ξ_i^P . The pattern statistics can thus be summarised as follows: Approximately 50% of units receive positively coded inputs, with the remaining units receiving no input for each intensity (Figure S1A, displaying 50 sample profiles $\vec{\xi}_i$ with different colors). The average input pattern activity exhibits a slight decay towards $\xi^{P/2}$ before increasing again (Figure S1B). The fraction of units receiving a non-zero input remains roughly constant across patterns (Figure S1C). The correlation structure between patterns exhibits a monotonic shape, which explains why the dynamics observed in the network after progressive Hebbian learning can be likened to a bell shape when considering the overlap with the inputs. In Figure S1D, we show, in a sample with $P = 50$, the overlap $\vec{\xi}^{20}$ (light gray), $\vec{\xi}^{25}$ (black), and $\vec{\xi}^{30}$ (dark gray) with all other patterns.

Simulations

All simulations were conducted in Python, and we provide code to run them online. We employed the following hyperparameters for eqs. (1) and (2): $\gamma = 0.1$, $\lambda = 0.014$, $\eta = 2$, and $g = 1.2$, $\kappa = N$, facilitating a homeostatic regime where learning and dynamic evolution proceed continuously without stagnation. For simplicity, we set the limiting average activity a , towards which the network's average activity is driven by the threshold, to match the average input activity. This value depends on the random realization of the input patterns in each simulation and is given by $a = \frac{1}{NP} \sum_{i=1}^N \sum_{\mu=1}^P \xi_i^\mu$, typically yielding approximately 0.2. Variations of these parameters are also feasible. We simulated an RNN with $N = 200$ units and parametrise the spectrum of weak to strong vibrational stimuli with $P = 100$ input patterns. We initialised the RNN weight matrix with random weights ranging from -1 to $+1$ at the beginning of each "experimental session", i.e. at the beginning of each run. Averages were performed on 100 random realisations each composed by 500 trials for the working memory paradigm, 200 random realisations each composed by 300 trials for the reference memory paradigm and 150 random realisations each composed of 750 trials for the one-back task. For the reproduction of the working memory paradigm of Akrami et al.³ we use $N = 204$, $P = 102$ in Figures 7A and 7B and $N = 196$, $P = 98$ in Figures 7C and 7D). Averages were performed over 50 random realisations each composed of 500 trials. The temporal parameters of the experimental protocol are configured as follows.

Working Memory Task 5 time steps before each stimulus, 5 time steps for stimulus presentation, followed by randomly selected intervals of 5, 10, or 20 time steps before the second 5-time step stimulus, followed by a randomise interval of 10-15 time steps before the following trial.

Reference Memory Task 10 time steps before the stimulus, 10 time steps for stimulus presentation, and a randomised interval of between 10 to 20 time steps after the stimulus.

One-Back Task 10-time steps for stimulus presentation, followed by randomly selected intervals of 20, 40, or 80-time steps from one stimulus to the next. In Figure 7 and we used a different stimulation protocol as reported in the caption.

Working Memory Task of Akrami et al.³ 5 time steps before each stimulus, followed by a presentation lasting 10 time steps for Figures 7A and 7B and 5 time steps for Figure 7D. Subsequent to the initial stimulus, random intervals of 15, 30, or 45 time steps occur before the second 10 steps stimulus, with a random interval between 10-15 time steps between subsequent trials.

Note that the exact relationship between ure cannot exceed one 8.5" x 11" page. Please do not include separate panels on multiple pages/files.

We ask that declaration these time steps and real milliseconds is beyond the scope of this article. However, we demonstrate in appendix further details on time scales in the model that this specific choice is not a fundamental feature to obtain the results shown in the article.

Reading out the network

We read out the intensity which is internally represented by the RNN model at each time as the intensity value corresponding to the stimulus that has the largest cosine similarity, or overlap, with the RNN state $\vec{V}(t)$. We perform the choices based on these representations as follows. In the working memory and in the one back paradigms, we compare the internal representation to the delivered stimulus at the stimulus onset (only for S_B in working memory). If the intensity of the delivered stimulus is larger than the internal representation of the network, the delivered stimulus is considered stronger than the previous one, and vice versa. In the reference memory paradigm, we consider the relaxation dynamics of the internal representation at the stimulus offset. If the internal representation relaxes towards higher intensities by the end of the trial, the stimulus is considered weak, and vice versa. When the choice is ambiguous, i.e. if the (second) stimulus coincides with the internal representation for working memory and one-back, or if the internal representation does not relax in reference memory, we take a random choice. Furthermore, we introduce "lapses" in all three paradigms, which reflect the animal's tendency to make errors due to a lack of concentration. Specifically, we take at random the choice in 10% of the trials, which are picked at random.

Fitting psychometric curves to data

We fitted the psychometric curves in the experimental and the simulated data of Figures 2, 4, and 7 using a the same modified sigmoidal function as Hachen et al.,¹⁰ which has the form

$$y = \gamma + (1 - \gamma - \lambda) \frac{1}{2} \operatorname{erf} \left(\frac{x - \mu}{\sqrt{2}\sigma} \right), \quad (\text{Equation 3})$$

where γ, λ, μ , and σ are all fitted to data (see Hachen et al.¹⁰ for a detailed discussion of the motivation and interpretation of the fitting parameters). In Figures 2 and 7, y is the probability that the second stimulus is judged stronger than the first, $S_B > S_A$, and x is a categorical variable indicating the difference between stimulus intensity. In Figure 4, y is the probability that a given stimulus is judged “strong”, and x is a categorical variable indicating the stimulus label with the weakest stimulus corresponding to 0.

Stimulation protocol

The stimulation protocol in the four different tasks dictates the stimuli delivered and the probabilities of delivering each of them at every step, as in the experiments. In the model, at every step in time, the input $\xi(t)$ appearing in Equation 1 is defined as follows:

$$\begin{cases} \xi(t) = \bar{\xi}^\mu & \text{if stimulus intensity labelled as } \mu \text{ is delivered} \\ \xi(t) = \bar{0} & \text{otherwise} \end{cases} \quad (\text{Equation 4})$$

The precise time-steps defining the pre-stimulus interval, inter-stimulus interval (ISI), and inter-trial interval (ITI) for each paradigm are provided in the STAR Methods section.

Stimulation setup stochasticity

For the tactile experiments, we also account for the additional variability introduced by the setup. While this is a minor detail, we include it to ensure a proper comparison between theory and experiment. As shown in the main text, the effect of this stochasticity is only mildly noticeable in the one-back paradigm, due to the very close-by intensities we selected.

In the experiments, as described in previous papers, see for example Hachen et al.,¹⁰ tactile inputs are delivered through a motor as a speed of vibration in a stochastic manner. A stimulus intensity is measured in terms of “average speed” \bar{v} . Each “average speed” results from a certain standard deviation σ . In fact, every 0.1 ms (10 kHz), a velocity drawn from a Gaussian distribution $|\mathcal{G}(0, \sigma)|$, where $|\cdot|$ is the absolute value, is delivered to the subject. The Gaussian distribution is illustrated for the set of stimuli delivered in the reference memory paradigm studied in Hachen et al.¹⁰ in Figure S1E. Negative values are converted to positive ones, resulting in a distribution with an effective mean $m(v) = \sqrt{2/\pi}\sigma$ (vertical line in Figure S1F). Through a motor-gain parameter y , the effective average velocity is then converted to $\bar{v} = ym(v) = 10\sigma$. Basically, in the experiments, any time a certain intensity is delivered to the network, this may be a slightly different average speed (i.e. intensity) from one time to the other. We reproduce, in Figure S1G, a sample distribution of the actual effective average velocities delivered across the different trials in a putative session of the published dataset in Hachen et al.¹⁰

In the model, we replicate the same bias in the distribution of the stimuli as follows: whenever an input defined by a standard deviation σ is delivered, we compute the average of 5000 random numbers drawn from the Gaussian distribution $|\mathcal{G}(0, \sigma)|$. This average give the effective input intensity which is delivered. We pick the corresponding discretised input pattern $\bar{\xi}^\mu$ accordingly. The quantity of 5000 was chosen to approximate the average intensity perceived over approximately 500ms.

In the figures we refer to the input intensity delivered (corresponding to the actual average velocity), and not to the standard deviations, in alignment with previous papers.

Synaptic Weight Changes and Biological Plausibility

In the main text, eq. (2) consists of two terms: a homeostatic term and a standard Hebbian learning term. The homeostatic term (the first one) is crucial for preventing the divergence of synaptic weights during the learning process, thus ensuring the stability of the network. As a result, the weights can change from positive to negative over time. Although synaptic weights in the model are intended to represent the effective connectivity between neurons—including potential contributions from inter-neurons that add complexity—we believe it is still important to address the phenomenon of weight switching, which is in apparent contradiction with Dale’s principle. In this section, we provide a more detailed discussion of this behaviour and present preliminary analyses to clarify the role of these weight changes within our model.

Figure S2 presents an analysis of synaptic weight changes across the three example tasks shown in Figure 6 of the main text. Our preliminary analysis reveals that approximately 20% of the total synaptic weights (as shown in Figure S2A) undergo at least one switch after a full trial of any paradigm. Interestingly, the maximum percentage of synaptic connections that change from one step to the next is 0.0125% ($100 \times 5 / 200^2$), as illustrated in Figure S2B. These weight switches are predominantly observed in synapses with low absolute weight values (not shown here), which are more likely to surpass the zero boundary during the homeostatic regulation process.

Recent studies have provided evidence for the presence of multi-transmitter neurons in the mammalian brain. For example, the work by Wallace and Sabatini²⁸ points to the existence of neurons that can switch between excitatory and inhibitory effects based on the relative presence of GABAergic and glutamatergic post synaptic receptors and transmitter received. This phenomenon, initially observed in the 1990s,^{42,43} and also more recently explored by Morales and Margolis,⁴⁴ suggests that some neurons may exhibit both excitatory and inhibitory activity under different circumstances. Nevertheless, introducing Dale’s law, which posits

that each neuron typically releases only one type of neurotransmitter, in our model, is an interesting direction for future work, offering deeper insights into the dynamic behaviour of synaptic weights.

Increasing synaptic homeostasis reproduces the effect of optogenetic inactivation of PPC³

We reproduce the experimental results of Akrami et al., obtained in a working memory paradigm, by utilising the same model as described in the main text, including the same parameters, but with stimuli and inter-stimulus interval distributions similar to those presented in Akrami et al.³ Our simulations reproduce the observed contraction bias effects (Figure 7A), in particular we show the probability of correct choices for 8 stimulus pairs (S_A , S_B) and the variations in the probability of correct choice among different pairs illustrate the presence of contraction bias, aligning with results depicted in Figure 1B of Akrami et al.³

We also reproduce the repulsive effect of previous stimulus pairs (Figure 7B), in particular we show how the probability of correct choice is influenced by the stimulus pair presented in the previous trial, showcasing a trend where lower intensity pairs result in a higher probability of S_B being chosen stronger than S_A . This mirrors the findings in Figures 2A–2C of Akrami et al.³

Furthermore, we reproduce the dependence of contraction bias on inter-stimulus intervals (Figure 7C), in particular we show the performance across a subset of 6 stimulus pairs, each with a constant $S_B = 76$ and S_A ranging from 66 to 84. Results demonstrate the impact of ISI intervals on performance, with shorter intervals leading to better performance, consistent with Figure 1E of Akrami et al.³

Finally, Akrami et al.³ showed that inactivating the posterior parietal cortex (PPC) leads to a reduction in contraction bias while maintaining memory capability. In our model, which normally operates in a homeostatic regime where learning and transient neural evolution occur continuously, we observe a similar effect if we modify the synaptic homeostatic parameter λ in Equation 2. This adjustment leads to a regime where synapses quickly forget previous stimuli. As a result, the internal representation in the network tends to become stuck to the last perceived stimulus, which effectively becomes the new fixed point on average. The attractor is thus very close to the last perceived stimulus and the drift of the memory trace during the ISI is minimal, consequently diminishing contraction bias while preserving memory until the next stimulus, see Figure 7D. In particular, we show how performance varies with changes in the synaptic homeostasis parameter λ . Increasing λ induces higher performance by accelerating synaptic forgetting of previous stimuli and facilitating faster attractor updates. This pathological regime replicates effects observed during the optogenetic inactivation of PPC, akin to Figure 3C of Akrami et al.³

Asymmetric shift of psychometric curves in Hachen et al.¹⁰ and individual differences

Figure 4B in the main text reproduces a plot of Hachen et al.¹⁰ The reported shift in the psychometric curves is asymmetric: the psychometric curves obtained when more high stimuli were delivered during the session (dark blue) has a point of subjective equality (square) corresponding to the new average; instead, the point of subjective equality of the psychometric curves obtained from these sessions where more low stimuli were delivered (light blue), do not coincide to the new average and is actually even slightly shifted to the right as compared to the PSE for uniform distributions. On the contrary, the result of the model is a symmetric shift, as reported in Figure 4F. We asked further details to the authors of Hachen et al.¹⁰ First they reinforced that the major result reported by their figure is the effective separation between the psychometric curves, which we reproduce, more than the quantitative distance of this shift. Second, they reported that the seemingly absence of movement of the light blue psychometric curve, with respect to an hypothetical uniform distribution of stimuli, is attributed to the tendency of the tested subjects of underestimation, i.e. of psychometric curves with PSE shifted towards the right. Thus, it remain open the question on whether a larger subject pool would actually lead to a symmetric shift in the psychometric curves as that we predict through the model. Note that in the results of the model we use a large sample size, thus obtaining smooth results, but if we would focus only on a subset of seeds which give rise, by chance, to right-shifted psychometric curves, we would also reproduce an asymmetric shift. We can see individual differences in our model as different realisations of the model dynamics, which will differ from each other due to different random settings (see STAR Methods). Randomness arise principally from the initial weights (akin to different subjects) and different random realisations of the input patterns and stimuli. In Figure S3, we report four examples of the psychometric curves and GLM estimates for four different realisations, composing, together with other 196 realisations, the average values presented in the plots in Figure 4E–4H. Individual peculiarities and stimulus history may lead to a different effects in the psychometric curves: for example in the realisation Figure S3A, there is a bias towards a strong categorisation, i.e. the attractor had a tendency to get stuck at high intensity, as opposed to the sample in Figure S3B.

Finally, we hypothesise that if experimentally, even with a larger pool, the shift would keep being asymmetric another component of the model which could contribute in reproducing the behaviour is the distribution of the input pattern profiles. Indeed, in all results we presented so far, a ratio of 50% is considered between positively and negatively coding stimuli. Changing this percentage may have some effect of this type on the learned patterns.

The effect of the stochasticity in the stimulation setup

We note that the asymmetry in the results for stimuli from the high vs. the low cloud in Figure 5D, as well as the decay of the red curve to a value larger than 0.5 in Figure 5A, arise from two independent factors. Firstly, the Weber-Fechner law⁴⁵ which suggests that the subjective sensation of stimulus intensity is proportional to its logarithm. This phenomenon facilitates the comparison of two equally distanced quantities when they are small compared to when they are large. Secondly, the stochasticity inherent in the stimulation setup which generates a broader distribution of inputs when stimuli are high, cf. Figure S1. In the model, but not in the experiments, we have the capability to eliminate the stochasticity inherent in the inputs, exposing the network solely to the average stimuli, thereby

rectifying the asymmetries in the results, as shown in [Figures S4A–S4C](#). Note that this effect is observed only in the one-back task because, in this design, the difference between consecutive intensities is much smaller than in the other experiments (see [STAR Methods](#)).

Further details on time scales in the model

The ratio between inter-stimulus interval and stimulus duration

In the results presented in the main text, the ratio between stimulus duration and inter-stimulus interval (ISI) in our simulations was not consistent with the experiments. Here we show that our model also reproduces the contractive bias in the working memory task even when the ratio between stimulus duration and ISI matches the experiment. [Figures S4F–S4J](#) reproduces the results from our simulation of the working memory task for convenience, as reported in the main paper. [Figures S4K–S4O](#), instead, reports the results of a simulation where we set the time step equal 333.5ms, resulting in 1 time step before each stimulus, 1 time step for stimulus presentation, followed by randomly selected intervals of 6, 12, or 24 time steps before the second 1-time step stimulus. All these durations precisely match the ratio of stimulus duration and ISI used in the experiments on real human subjects. Since the stimulus is applied for a much shorter time, we increased the stimulus intensity to $\eta = 5$ compared to the simulations reported in the main text. In the main text, we report results from simulations with ratios that do not closely match experiments for computational efficiency: applying the stimulus for longer, with a smaller η , allowed us to perform trials in less total timesteps while having a stimulus duration that was long enough to track the dynamics of its effect on the attractor. We made this choice based on experimental evidence from Toso et al.,⁴⁶ who showed that the percept of a stimulus is independent of its duration for a stimulation length above 300 to 500ms, see [Figure 5D](#) in Toso et al.⁴⁶ We therefore assumed that we could extend the stimulus duration in our simulations without changing how the modelled subject would have experienced its intensity. As we show in the next subsection, our model reproduces other behavioural effects related to timescales with the hyperparameter settings used in the main text.

The RNN model qualitatively reproduces behavioural effects of varying ISI

Experimental evidence shows that contraction bias depends on the time between stimuli – the longer the interval, the stronger the contraction bias. For instance, Akrami et al.³ found that rats presented with pairs of auditory stimuli at varying ISIs (2, 4, 6 seconds) showed a characteristic time dependence in their psychometric curves: longer ISIs lead to less steep curves, as was reported in [Figure 1E](#) of their paper. This behaviour is captured in our model, as shown in [Figure 7C](#).

We also found that our model correctly reproduces the behavioural effect of varying the inter-trial interval (ITI) in the One-Back paradigm with tactile stimuli, where the ITI is equal to the inter-stimulus interval. Despite differences in species and sensory modality, conditioned psychometric curves show the typical contraction bias effect, as shown in [Figure S4D](#), as the ITI is varied from 2 to 4 to 8. Our model, using proportionally varying ITIs (20, 40, 80 time steps), replicates this behaviour, as shown in [Figure S4E](#).

Neural vs. memory time scales

In our model, threshold linear units are a simplified representation of neuronal responses and are not directly mapped to intrinsic neuronal properties. They can represent individual neurons or the ensemble rate of a group of interconnected neurons. If we interpret the threshold linear units as modeling individual neurons, and considering the relationship between γ and synaptic conductances, then the time steps used in our simulations (e.g., 335 ms) would imply unrealistically large synaptic conductances on the order of seconds. One way to address this would be to either increase γ (which destabilises the simple neural model) or increase the number of time steps so that each corresponds to a duration between 1–4 ms, leading to a more realistic time scale of synaptic conductance. This would require a retuning of the model's hyperparameters, for which a better theoretical understanding of the different regimes of our model would be of great help. However, this is out of the scope of the present work. Here, we show that even with the present parameters, our model addresses the key challenge of modelling working memory, which is overcoming the disparity between neural time scales and memory duration. In our model, the duration for which a stimulus is held in memory and the rate at which its memory trace diverges depend not only on the neuronal time scale but also on the attractor's shape and its ongoing modification through Hebbian learning solving the highlighted challenge. In [Figure S5A](#), we illustrate how the memory trace of different stimuli S_A evolves from their offset to the step just before the onset of S_B , across all trials with the longest ISI (24 time steps) for a single model realization. The variation in the internal representation of the stimulus across trials as well as how fast they diverge results from both the activity state prior to stimulation and the strength of the recurrent connections. The vertical line at 10 time steps indicates the critical value that would be expected if memory were solely dependent on the neuronal time scale (assuming $\gamma = 0.1$), but no significant change occurs at this point. In [Figure S5B](#), where we average across trials and realizations, this effect becomes even more apparent. The memory of the stimulus persists well beyond the neuronal critical value, as the curves do not collapse either at 10 time steps or at 24.

While neuronal time scales smooth the transition of activity and help maintain a stimulus representation, the free evolution of the system, in the absence of stimulation, is governed by the attractor, or the connections, which store the history of previous inputs. The basin of attraction can take different shapes—such as the steeper or smoother parabolas one could draw in [Figures 3E and 3F](#) of the main text—which continuously change depending on the inputs and the system's dynamics. A steeper basin causes the memory to decay more quickly, while a smoother one slows this process. Additionally, depending on the stimulus history, the representation of a stimulus may or may not align with the attractor in a given trial, leading to different memory divergence behaviours, despite the same neuronal time constant. Therefore, the memory time scale cannot be directly mapped to the neuronal time scale; it is instead an **effective time scale** that emerges from these factors and changes throughout the experiment.