

Synaptic mechanisms of interference in working memory

Zachary P. Kilpatrick^{1,2}

Information from preceding trials of cognitive tasks can bias performance in the current trial, a phenomenon referred to as interference. Subjects performing visual working memory tasks exhibit interference in their responses: the recalled target location is biased in the direction of the target presented on the previous trial. We present modeling work that develops a probabilistic inference model of this history-dependent bias, and links our probabilistic model to computations of a recurrent network wherein short-term facilitation accounts for the observed bias. Network connectivity is reshaped dynamically during each trial, generating predictions from prior trial observations. Applying timescale separation methods, we obtain a low-dimensional description of the trial-to-trial bias based on the history of target locations. Furthermore, we demonstrate task protocols for which our model with facilitation performs better than a model with static connectivity: repetitively presented targets are better retained in working memory than targets drawn from uncorrelated sequences.

Parametric working memory experiments are a testbed for behavioral biases and errors, and help identify neural mechanisms that underlie them¹⁻³. In visuospatial working memory, subjects identify, store, and recall target locations in trials lasting a few seconds. Response errors are normally distributed⁴⁻⁶, and tend to accumulate during the delay-period, while subjects retain the target location in memory^{1,6,7}. Complementary neural recordings suggest these working memories are implemented in circuits comprised of stimulus-tuned neurons with slow excitation and broad inhibition^{8,9}. Persistent activity emerges as a tuned pattern of activity called a bump state, whose peak encodes the remembered target position^{6,10}.

Neuronal studies of visual working memory typically focus on population activity within a single trial, ignoring serial correlations across trials¹¹. Several authors have identified behavioral biases that cause the previous trial's visual target to interfere with the subject's response on the subsequent trial^{12,13}. For instance, in delayed match-to-sample tests, false alarms occur more often when comparison stimuli match samples from previous trials¹⁴.

As has been shown previously, Eq. (2) can be written iteratively³⁶:

$$L_{n+1} = \frac{P(1:n-1)}{P(1:n)} f_n(\cdot) L_n,$$

suggesting such a computation could be implemented and represented by neural circuits. Temporal integration of tuned inputs has been demonstrated in both neural recordings³⁷⁻³⁹ and circuit models

Note, sequential computations are trivial in the limit of a constantly-changing environment¹, since the observer assumes the environment is reset after each trial. Prior observations provide no information about the present distribution, so the predictive distribution is always uniform: $L_{n+1} = \bar{P}_0$.

In summary, a probabilistic inference model that assumes the distribution of targets is predictable over short timescales leads to response biases that depend mostly on the previous trial. We now demonstrate that this predictive distribution can be incorporated into a low-dimensional attractor model which describes the degradation of target memory during the delay-period of visual working memory tasks^{10,41,42}.

Incorporating suboptimal predictions into working memory. We model the loading, storage, and recall of a target angle using a low-dimensional attractor model spanning the space of possible target angles $[-180, 180]^\circ$. These dynamics can be implemented in recurrent neuronal networks with slow excitation and broad inhibition^{6,9,43}. Before examining the effects of neural architecture, we discuss how to incorporate the predictive distribution update, Eq. (3), into an associated low-dimensional model. Our analysis links the update of the predictive distribution to the spatial organization of attractors in a network. Importantly, working memory is degraded by dynamic fluctuations, so the stored target angle wanders diffusively during the delay-period^{6,9,42}.

During the delay-period of a single trial, the stored target angle $\theta(t)$ evolves according to a stochastic differential equation¹⁰:

$$d\theta(t) = -\frac{d\langle \theta(t) \rangle}{d\theta} dt + \sigma dW(t). \quad (4)$$

Here $\langle \theta(t) \rangle$

$q(x, t)$ determines an evolving potential function

at $n+1$. The STF variable's center-of-mass $q(t)$ slowly drifts towards m which allows (t) to drift there as well, $\overline{(- q(t))}$. This accounts for the slow build-up of the bias that increases with the length of the delay-period¹³.

distinct from neural activity¹³, as dynamic synapses are in our model. In total, our model provides both an intuition for the behavioral motivation as well as neurophysiological mechanisms that produce such interference.

Discussion

Comparison with previous work. The work of Papadimitriou *et al.*^{13,55} also contains modeling studies, accounting for some aspects of their experimental observations. Our computational model differs from and extends their findings in several important ways. We propose that interference can arise as a suboptimal inference

Methods

Assumptions of the inference model.

is population rate model can be explicitly analyzed to link the architecture of the network to a low-dimensional description of the dynamics of a bump attractor as described by Eq. (4).

Each location x in the network receives recurrent coupling defined by the weight function $w(x - y)$ via a convolution $w(x) \cdot g(x) = \int_{-180}^{180} w(x - y)g(y) dy$. We take this function to be peaked when $x = y$ and decreasing as the distance $|x - y|$

in¹³. Intertrial intervals are varied to produce Fig. 5B by drawing $T_I^n := t_{n+1} - (T_C + T_D^n + T_A)$ randomly from a uniform pmf for the discrete set of times $T_I^n \in \{1000, 1200, \dots, 5000\}$ ms and θ_n randomly as in Fig. 5A and identifying the θ_n that produces the maximal bias for each value of T_I^n . Delay-periods are varied to produce Fig. 5C by drawing T_D^n randomly from a uniform pmf for the discrete set of times $T_D^n \in \{0, 200, \dots, 5000\}$ ms and following a similar procedure to Fig. 5B. Draws from a uniform density function $P(\theta_n) = \frac{1}{360}$, defined on $\theta_n \in [-180, 180]^\circ$ are used to generate the distribution in Fig. 6A and plots in Fig. 7. Nontrivial correlation structure in target selection is defined by the sum of a von Mises distribution and uniform distribution $\text{corr}(\theta_{n+1}, \theta_n) = (1 - \epsilon) e^{25 \cos(\theta_n - \theta_{n+1} - \mu)} + \epsilon P_0$ for fixed ϵ with $\epsilon = 0.5; \mu = 0$ for local correlations (Fig. 6B) and $\mu = 90$ for skewed correlations (Fig. 6C).

The recurrent network, Eq. (15), is assumed to encode the initial target θ_n during trial n via the center-of-mass $\langle \theta(t) \rangle$ of the corresponding bump attractor. Representation of the cue at the end of the trial is determined by performing a readout on the neural activity $u(x, t)$ at the end of the delay time for trial $n: t = t_n + T_C + T_D^n$. One way of doing this would be to compute a circular mean over x weighted by $u(x, t)$, but since $u(x, t)$ is a roughly symmetric and peaked function in x , computing $\langle \theta(t) \rangle := \text{argmax}_x u(x, t)$ (when $t \in [t_n, t_n + T_C + T_D^n)$) is an accurate and efficient approximation^{6,42}. The bias and relative saccade endpoint on each trial n are then determined by computing the difference $\langle \theta(t) \rangle - \theta_n$ (Figs 5, 6 and 7).

Deriving the low-dimensional description of bump motion. We analyze the mechanisms by which STF shapes the bias on subsequent trials by deriving a low-dimensional description for the motion of the bump position $\langle \theta(t) \rangle$. To begin, note that in the absence of facilitation ($\epsilon = 0$), the variable $q(x, t) = 0$. In the absence of noise ($W(x, t) = 0$), the resulting deterministic Eq. (15) has stationary bump solutions that are well studied and defined by the implicit equation^{43,47,89}:

$$U(x) = \frac{1}{2} \left(\frac{1}{\cos(x - \theta_n)} - \frac{1}{\cos(x - \theta_{n+1})} \right)$$

$$K(x, t_{n+1}) = \frac{df(x)}{dt}, \quad (19)$$

where K is a scaling constant and t_{n+1} is the starting time of trial $n + 1$ in the original time units $t = t_s / \omega$. The form of the probability $f(x)$ that can be represented is therefore restricted by the dynamics of the facilitation variable $q(x, t)$. We can perform a direct calculation to identify how $q(x, t)$ relates to the predictive distribution it represents in the following special case.

Explicit solutions for high-gain firing rate nonlinearities. To explicitly calculate solutions, we take the limit of high-gain, so that $F(u) \approx H(u - \theta)$ and $w(x) = \cos(\omega_1 x)$, note $\omega_1 = 180^\circ / \lambda$. Note, we have compared our predictions here with the results of numerical simulations for sigmoidal firing rates $F(u) = 1/[1 + e^{-(u-\theta)}]$ with gain $\omega_1 = 20$, and the results are in good agreement. In this case, the bump solution $U(x - x_0) = (2 \sin(a) / \omega_1) \cos(\omega_1(x - x_0))$ for $U(\pm a) = 1$ and null vector $V(x - x_0) = (x - x_0 - a) - (x - x_0 + a)$ (without loss of generality we take $x_0 = 0$)⁴⁷. Furthermore, we can determine the form of the evolution of $q(x, t)$ by studying the stationary solutions to Eq. (15) in the absence of noise ($W = 0$). For a bump $U(x)$ centered at $x_0 = 0$, the associated stationary form for $Q(x)$ assuming $H(U(x) - \theta) = 1$ for $x \in (-a, a)$ and zero otherwise is $Q(x) = q_+ / (1 + \omega_1)$ for $x \in (-a, a)$ and zero otherwise. Thus, if the previous target was at x_m , we expect $q(x, t)$ to have a shape resembling $Q(x - x_m)$ after trial n . Assuming the cue plus delay time during trial n was $T_C + T_D^n$ and the intertrial interval is T_I^n , slow dynamics will reshape the amplitude of $q(x, t)$ so $q_n(T^n) = (1 - e^{-(T_C + T_D^n)/\tau}) e^{-T_I^n/\tau}$ ($T^n = T_C + T_D^n + T_I^n$ is the total time g]TJ EMC /T6 tt

68. Lim, S. & Goldman, M. S. Balanced cortical microcircuitry for maintaining information in working memory. *Nat. Neurosci.* **16**, 1306–1314 (2013).
69. Boerlin, M., Machens, C. K. & Denève, S. Predictive coding of dynamical variables in balanced spiking networks. *PLoS Comput. Biol.* **9**, e1003258 (2013).
70. Shaham, N. & Burak, Y. Slow diffusive dynamics in a chaotic balanced neural network. *PLoS Comput. Biol.* **13**, e1005505 (2017).
71. Ma, W. J., Husain, M. & Bays, P. M. Changing concepts of working memory. *Nat. Neurosci.* **17**, 347–356 (2014).
72. Nassar, M. R., Helmers, J. C. & Frank, M. J. Chunking as a rational strategy for lossy data compression in visual working memory tasks. *bioRxiv* 098939 (2017).
73. Zhang, W. & Luck, S. J. Discrete fixed-resolution representations in visual working memory. *Nature* **453**, 233–235 (2008).
74. Luck, S. J. & Vogel, E. K. Visual working memory capacity: from psychophysics and neurobiology to individual differences. *Trends Cogn. Sci.* **17**, 391–400 (2013).
75. Bays, P. M. & Husain, M. Dynamic shifts of limited working memory resources in human vision. *Science* **321**, 851–854 (2008).
76. Wei, Z., Wang, X.-J. & Wang, D.-H. From distributed resources to limited slots in multiple-item working memory: a spiking network model with normalization. *J. Neurosci.* **32**, 11228–11240 (2012).
77. Almeida, R., Barbosa, J. & Compte, A. Neural circuit basis of visuo-spatial working memory precision: a computational and behavioral study. *J. Neurophysiol.* **114**, 1806–1818 (2015).
78. Bays, P. M. Spikes not slots: noise in neural populations limits working memory. *Trends Cogn. Sci.* **19**, 431–438 (2015).
- 79.