

# Synaptic Correlates of Working Memory Capacity

## Highlights

- Multiple items can be brought to working memory via periodic brief reactivations
- Working memory capacity can be analytically estimated in the framework of this model
- Capacity scales as time of synaptic depression over that of synaptic current
- Stream of inputs can be segmented into chunks by modulating external excitation

## Authors

Yuanyuan Mi, Mikhail Katkov,  
Misha Tsodyks

## Correspondence

misha@weizmann.ac.il

## In Brief

Mi, Katkov, and Tsodyks derived an approximate analytical expression for working memory capacity in the framework of synaptic theory. This development predicts how manipulating the parameters of short-term synaptic plasticity and synaptic time constant will affect the capacity.



# Synaptic Correlates of Working Memory Capacity

Yuanyuan Mi,<sup>1,2,3</sup> Mikhail Katkov,<sup>2</sup> and Misha Tsodyks<sup>2,3,4,\*</sup><sup>1</sup>Brain Science Center, Institute of Basic Medical Sciences, Beijing 100850, China<sup>2</sup>Department of Neurobiology, Weizmann Institute of Science, Rehovot 76100, Israel<sup>3</sup>Department of Neuroscience, Columbia University, New York, NY 10032, USA<sup>4</sup>Lead Contact\*Correspondence: [misha@weizmann.ac.il](mailto:misha@weizmann.ac.il)<http://dx.doi.org/10.1016/j.neuron.2016.12.004>

## SUMMARY

Psychological studies indicate that human ability to keep information in readily accessible working memory is limited to four items for most people. This extremely low capacity severely limits execution of many cognitive tasks, but its neuronal underpinnings remain unclear. Here we show that in the framework of synaptic theory of working memory, capacity can be analytically estimated to scale with characteristic time of short-term synaptic depression relative to synaptic current time constant. The number of items in working memory can be regulated by external excitation, enabling the system to be tuned to the desired load and to clear the working memory of currently held items to make room for new ones.

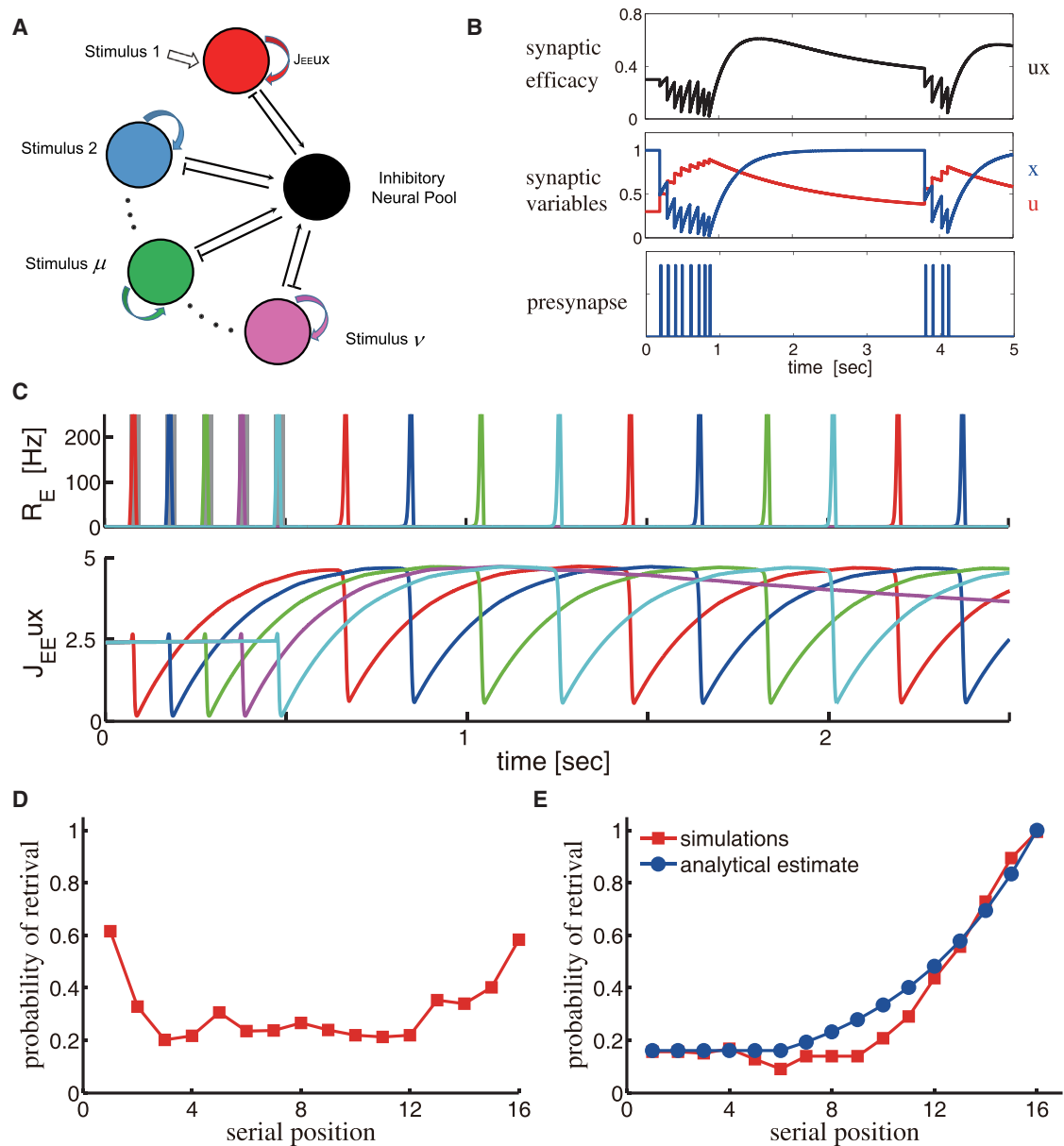
## INTRODUCTION

Working memory (WM) refers to short-term storage and manipulation of information (Miller et al., 1960; Baddeley and Hitch, 1974; Baddeley, 2003). There is hardly a cognitive task that does not involve WM, including visual processing, speech comprehension, and episodic memory (Cowan, 2001). Nevertheless, WM capacity is extremely limited, ranging between three and six items for most healthy human participants (Cowan, 2001; Fukuda et al., 2010; Luck and Vogel, 1997). It is often postulated that the brain possesses a specialized buffer, or “focus of attention,” where memory items can be temporarily placed for short periods of time and removed when needed; hence, WM capacity corresponds to the size of this buffer (Cowan, 2001; Oberauer, 2002). The neuronal implementation of the focus of attention and its size, as well as the way memory items can be placed and removed from it, are not understood. The most popular hypothesis is that WM is mediated by persistent activity of neurons encoding the corresponding items in long-term memory (see, e.g., Compte et al., 2000; Wei et al., 2012; Edin et al., 2009). The maximal number of items simultaneously active depends on the characteristics of the network in a complex way, but there does not seem to be a fundamental upper limit on WM capacity in this model (Amit et al., 2003; Rolls et al., 2013). Another view posits that WM involves sequential activations of item representations (Cowan, 2010; Horn and Opher, 1996; Raffone and Wolters, 2001; Lundqvist et al.,

2016; Lisman and Idiart, 1995). Lisman and Idiart (1995) suggested that WM involves periodic reactivation of memory representations at each gamma cycle within gamma-theta nested oscillations in hippocampus, mediated by slow after depolarization with time constant that should be matched to theta period. Assuming each memory is activated exactly once during a theta cycle and that the same memories are repeatedly reactivated over subsequent theta cycles (Siegel et al., 2009), the WM capacity is then estimated as a ratio of gamma and theta frequencies, which is compatible with earlier psychophysical estimates (Miller, 1956). The explicit dependence of WM capacity on the parameters of the model was not considered in this study, but it appears that the crucial factor limiting the capacity is the theta rhythm period, which could depend on several intrinsic cellular properties, e.g., nonselective cation channels (Colgin, 2013).

Recently, Mongillo et al. proposed a synaptic-based theory for short-term information storage in neural circuits (Mongillo et al., 2008; see also Lundqvist et al., 2011). In this model, memory is retained by item-specific pattern of synaptic facilitation. This mechanism does not require neurons to fire with elevated rate for the whole duration of the memory task, resulting in a robust and metabolically more efficient scheme. Several items can be maintained in the WM via consecutive brief reactivations of the corresponding neuronal groups. Here we aim to analytically estimate the maximal number of items that can be maintained in WM. The advantage of analytical expression is that it allows one to make predictions about how WM capacity depends on the various synaptic, neuronal, and circuit parameters that can potentially be tested by genetic manipulations.

A basic assumption of our model is that only one memory representation can be active at any single moment, which is guaranteed by strong reciprocal connections to a global non-specific inhibitory pool (see Figure 1A below), consistent with experimental data (Fino and Yuste, 2011). If each memory representation had its own inhibition and hence was independent of the others, there would be no fundamental constraint to the capacity beyond the overlaps between the representations. We could think of several reasons against this idea: (1) the brain should then have an additional processing step that would disentangle the simultaneously active populations into a set of memories, and (2) each time a new memory is stored in the network, inhibition should adapt to this by generating a new population of inhibitory neurons that are specific to a new ensemble. This would make it much harder to add new representations, i.e., learning would suffer to improve WM capacity.



**Figure 1. STP-Based WM Network Model**

(A) Network architecture: a number of recurrent excitatory neural clusters, shown in different colors, reciprocally connected to an inhibitory neuron pool, shown in black.

(B) Model of a synaptic connection with STP. In response to a presynaptic spike train (lower panel), the neurotransmitter release probability  $u$  increases and the fraction of available neurotransmitter  $x$  decreases (middle panel), representing, synaptic facilitation and depression, respectively. The effective synaptic efficacy is proportional to  $ux$  (upper panel).

(C) Network simulation with five loaded memory items. Upper panel: firing rates of different clusters. Five clusters are sequentially stimulated by brief external excitation (shaded colored rectangles). Different colors correspond to different clusters as in (A). Following the stimulation, four clusters continue sequential activation in the form of PSs while the remaining item fades away. Lower panel: the instantaneous synaptic efficacy  $J_{EEUX}$  for stimulated neural clusters during loading and subsequent reactivations.

(D) Probability of retaining an item in WM as a function of its serial position in a train of 16 loaded items, computed with 450 simulated trials with presentation frequencies between 16 and 66 items/s.

(E) Red: same as (D); computed with 450 simulated trials with presentation frequencies between 0.25 and 2 items/s. Blue: analytical calculation of the probability, as explained in the text and [Supplemental Information](#). The parameters are as follows:  $J_{EE} = 8$ ,  $\tau_f = 1.5s$ ,  $\tau_d = 0.3s$ ,  $U = 0.3$ ,  $\tau = 8ms$ ,  $J_{IE} = 1.75$ ,  $J_{EI} = 1.1$ ,  $\alpha = 1.5$ ,  $P = 16$ , and  $I_b = 3.0Hz$  in (C) and  $I_b = 8Hz$  in (D) and (E).

## RESULTS

We consider a neural network model with memory items encoded by interconnected neuronal ensembles with short-term plasticity (STP) of recurrent connections (Mongillo et al., 2008). To achieve an analytical hold on capacity estimation, we drastically reduce the complexity of the network to leave only the most essential features that allow it to function as WM. We neglect the overlaps between different representations so each item is represented by a single excitatory unit (cluster) characterized by its activity rate, with self-excitation reflecting the strengthened connections between neurons encoding a given item in long-term memory.

Following the STP model developed in Markram et al. (1998), recurrent excitatory connections are characterized by fixed “absolute synaptic efficacy” and two dynamic variables:  $u$ , which stands for release probability, and  $x$ , the fraction of available neurotransmitters (Figure 1B). If  $J_{EE}$  is the absolute synaptic efficacy between two excitatory neurons, the instantaneous synaptic efficacy subject to STP is given by  $J_{EE}ux$ . Upon arrival of a spike, the release probability  $u$  temporarily increases, resulting in short-term facilitation. Meanwhile, the fraction of available neurotransmitters  $x$  decreases, resulting in short-term depression. After neuronal spiking,  $u$  returns to its baseline value  $U$  with a time constant  $\tau_f$ , and  $x$  recovers to its maximum value  $x = 1$  with a time constant  $\tau_d$ .

In (Tsodyks et al., 1998), the STP model was used to derive the expression for the postsynaptic current resulting from the activity of a large, uncorrelated pre-synaptic population. The resulting network model has three differential equations for each of  $P$  excitatory clusters (synaptic current  $h_\mu$  and two STP variables  $u_\mu$  and  $x_\mu$  for each cluster  $\mu$ ;  $\mu = 1, \dots, P$ ) and one additional equation for the inhibitory pool current  $h_I$ :

$$\tau \frac{dh_\mu}{dt} = -h_\mu + J_{EE}u_\mu x_\mu R_\mu - J_{EI}R_I + I_b + I_e(t), \quad (\text{Equation 1})$$

$$\frac{du_\mu}{dt} = \frac{U - u_\mu}{\tau_f} + U(1 - u_\mu)R_\mu, \quad (\text{Equation 2})$$

$$\frac{dx_\mu}{dt} = \frac{1 - x_\mu}{\tau_d} - u_\mu x_\mu R_\mu, \quad \text{and} \quad (\text{Equation 3})$$

$$\tau \frac{dh_I}{dt} = -h_I + J_{IE} \sum_\nu R_\nu, \quad (\text{Equation 4})$$

where  $\tau$  is the neuronal time constant, for simplicity the same for excitation and inhibition;  $I_b$  is the constant background excitation that we assume to reflect the attentional state of the network (Zhang et al., 2014); and  $I_e$  is the external input used to load memory items into the network. As in Mongillo et al. (2008), we consider facilitating synapses with  $\tau_f \gg \tau_d$ .  $R(h) = \alpha \ln(1 + \exp(h/\alpha))$  is neuronal gain chosen in the form of a smoothed threshold-linear function, also the same for excitatory and inhibitory neurons. The exact shape of the gain function is not important, but it should exhibit a tail for negative currents in order for the network to generate population spikes (see below; Tsodyks, 2004). The tail in the gain function could

emerge either from noisy input that can cause some degree of firing for subthreshold levels of current, or when considering the effects of non-homogeneities in firing thresholds in neuronal clusters representing memory items.

To illustrate the proposed mechanism, we simulated the network with parameters that are compatible with experimental measurements of inter-pyramidal connections in the prefrontal cortex (Wang et al., 2006). We loaded five items into WM by applying transient external excitation inputs to the corresponding units and observed that four of them were maintained successfully in the form of brief reactivations called population spikes (PSs; Tsodyks et al., 2000), indicating that for this set of parameters, the capacity of WM is four (see Figure 1C). To further characterize the model, we subjected it to longer trains of stimuli at different presentation frequencies and computed the “retention curve,” i.e., the probability that an item with a given serial position in the train is retained in WM at the end of the train. We observed a non-monotonic retention curve for high presentation frequencies: first and last items had a higher chance of being retained (primacy and recency effect, respectively; Figure 1D); for low frequencies, there was a pronounced recency effect that could be estimated accurately by assuming that each time a new item is presented, one of the previous items is erased from WM with equal probability (Figure 1E; see Supplemental Information, available online, for details).

As seen in Figure 1C, WM containing several items corresponds to the periodic activity state of the network (“limit cycle”) where respective clusters emit PSs in the fixed order. As a dynamical system, the network exhibits multistability with many stable limit cycle solutions that evolve from different initial conditions. Rigorous analysis of coexisting limit cycle solutions is mathematically intractable, but this view allows us to address the issue of capacity numerically by simulating the network with random initializations until it converges to one of the periodic solutions (see Supplemental Information). In Table 1, we show the results of these simulations for different values of background input, with STP parameters that are compatible with experimental data from the prefrontal cortex (Wang et al., 2006). The table shows probabilities  $P_i$  that the network converges to a limit cycle with  $i$  sequentially activated clusters, computed with 200,000 random initializations of the network state at each background input value.  $P_0$  corresponds to a baseline state with no PSs. The results show that (1) only when the background input is high enough can multiple items be kept in WM, and (2) with an increasing background input, WM capacity is gradually increasing. For network parameters chosen in these simulations, the maximal capacity is six, since further increase in background input leads to destabilization of the baseline spontaneous state of the network.

Even with the extreme simplifications described above, the model is still described by a relatively large number of parameters, namely STP synaptic parameters, neuronal gain functions, synaptic strengths, and time constant; hence, “brute force” numerical approach would be impractical and unrevealing. We now present the intuitive outline of the derivation of the final result (see Supplemental Information for more details).

The maximum number of items that can be maintained in WM is determined by the ratio of two factors: (1) the maximal period

**Table 1. The Chances of Converging to a State with a Different Number of Items in WM for Different Arousal Levels**

$I_b$	$P_0$	$P_1$	$P_2$	$P_3$	$P_4$	$P_5$	$P_6$	$P_7$
2.4	1	0	0	0	0	0	0	0
2.45	0.9998	0.0002	0	0	0	0	0	0
2.5	0.9991	0.0008	0.0001	0	0	0	0	0
2.56	0.9950	0.0039	0.0010	0.0001	0	0	0	0
3.0	0.5668	0.1129	0.1420	0.1772	0.0011	0	0	0
3.7	0.0026	0.0151	0.0701	0.2872	0.6238	0.0012	0	0
5.5	0.0003	0.0008	0.0013	0.0242	0.2351	0.7015	0.0368	0
7	0.0002	0.0007	0.0008	0.0066	0.1187	0.6906	0.1824	0
14	0.0001	0.0004	0.0006	0.001	0.1387	0.8506	0.0086	0

The other parameters as in Figure 1.

$T_{\max}$  of the limit cycle of the network, i.e., the maximal time between subsequent reactivation of each cluster, and (2) the temporal separation between two consecutive PSs, referred to as  $t_s$ ; the capacity of WM is given by the maximum number of PSs that can be accommodated in a single period of the limit cycle, i.e., by

$$N_C \approx T_{\max}/t_s. \quad (\text{Equation 5})$$

To estimate  $T_{\max}$ , we note that a PS is triggered by intrinsic instability due to recurrent excitation; hence, it is always produced by the cluster that has, at that moment, the largest effective recurrent strength  $J_{EELUX}$ . Since the time evolutions of the effective strengths for different clusters have identical shape triggered by corresponding PSs, and separated from each other by the time  $t_s$ , which is significantly shorter than the  $T_{\max}$ , we conclude that the longest time between activations of a cluster approximately equals the time it takes for the synaptic efficacy curve to reach a peak (see Figures 1C and S2). This is determined by the solution of Equations 2 and 3 above, which can be greatly simplified by neglecting the firing rate of a cluster between the PSs, turning them into linear equations, and resulting in the following expression:

$$T_{\max} \approx \tau_d \ln \frac{\tau_f/\tau_d}{1-U}, \quad (\text{Equation 6})$$

i.e.,  $T_{\max}$  is chiefly determined by the time constant of synaptic depression, depending weakly on other STP parameters.

The  $t_s$  has three components: the width of the PS of the previous item, the delay and the width of the inhibitory pulse triggered by this PS, and finally, the time it takes for a next cluster to recover from inhibition and initiate the new PS (Figures 2A and S3). It can be intuited that the first two components are proportional to the synaptic time constant  $\tau$  (see Supplemental Information), and the third, dominant component should be found by solving Equations 1, 2, 3, and 4 above. To simplify these equations, we note that on the time-scale of  $t_s$ , which is significantly shorter than  $T_{\max}$ , we can neglect the STP evolution of effective synaptic strength for a cluster that is about to emit a PS (Equations 2 and 3) and replace it with its maximal value,  $J_{\max}$ . If we also neglect the time evolution of inhibition between the PSs ( $I_{\text{inh}} = J_{EI}R_i$ ),

we are left with the single dynamical equation for synaptic current of a cluster:

$$\tau \frac{dh}{dt} = F(h), \quad (\text{Equation 7})$$

$$F(h) = -h + J_{\max}R(h) + (I_b - I_{\text{inh}}). \quad (\text{Equation 8})$$

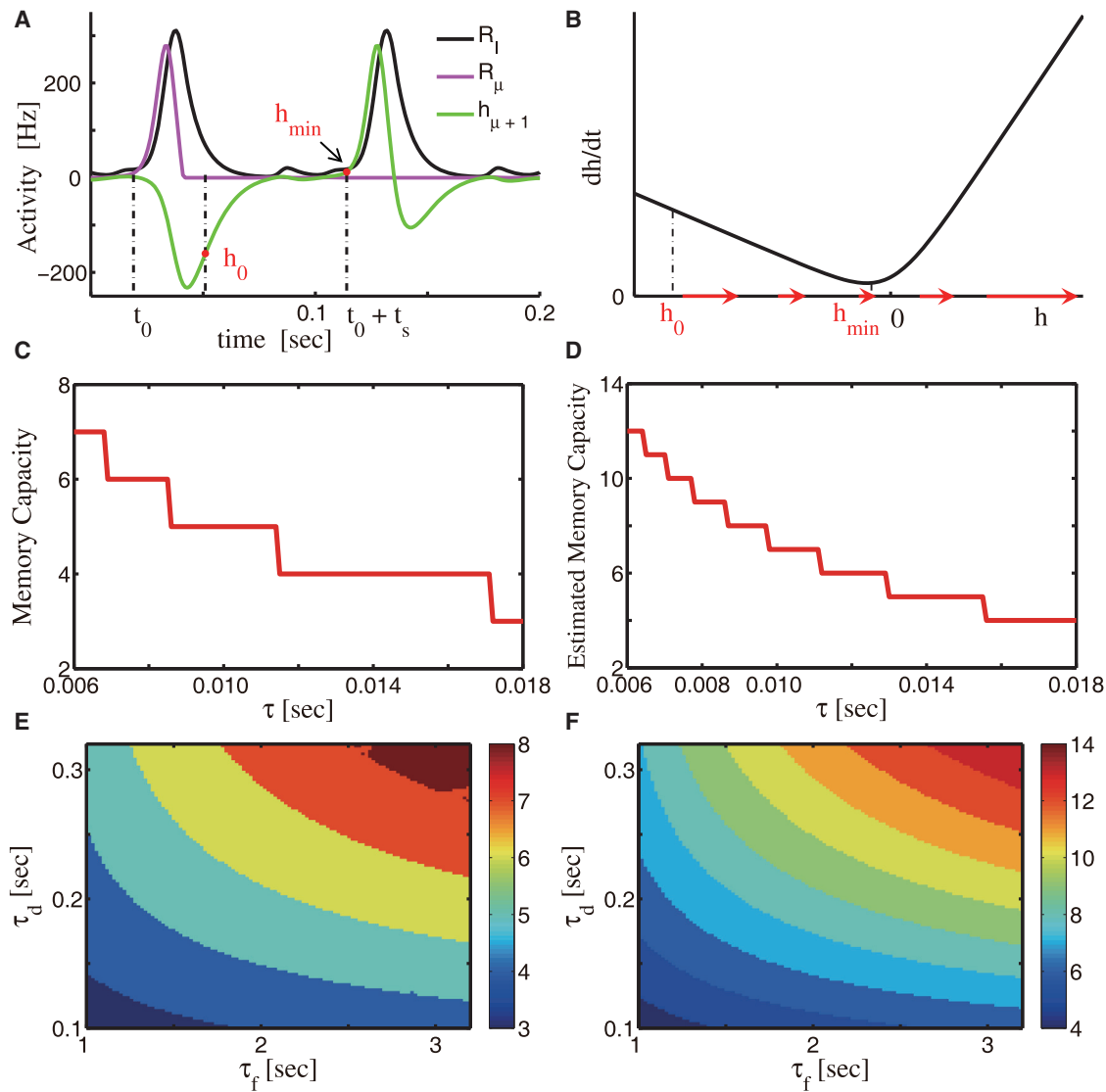
The function  $F(h)$  that determines the flow of synaptic current is composed of two approximately linear branches with negative and positive slopes, respectively (Figure 2B). The strongly negative initial value  $h_0$  is determined by the strength of inhibition triggered by the PS of the previously activated item. The time flow of the current due to Equation 7 will consist of two distinct segments: progressively slow recovery from inhibition until the minimum  $h_{\min}$  of  $F(h)$  is reached, followed by fast acceleration signaling the onset of the PS (see Figure 2B, red arrows). The time it takes for the cluster to emit the PS is thus given by the time for  $h$  to reach  $h_{\min}$ . Using the linear approximation for the negative branch, the analytical estimate for  $t_s$  can be computed as

$$t_s \approx \tau \left( \ln \frac{|h_0|}{I_b - I_{\text{crit}}} + C \right), \quad (\text{Equation 9})$$

where  $C$  reflects the contribution of PS width and inhibition duration, while  $I_{\text{crit}} \approx I_{\text{inh}} - \alpha \ln(J_{\max} - 1)$  is the critical value for the background excitation for which  $t_s$  diverges logarithmically (see Supplemental Information). Combining the above analysis results in the following analytical estimate for WM capacity:

$$N_C \approx \frac{\tau_d}{\tau} \frac{\ln \frac{\tau_f/\tau_d}{1-U}}{\ln \frac{|h_0|}{I_b - I_{\text{crit}}} + C}. \quad (\text{Equation 10})$$

Two conclusions can be drawn from this result: (1) the WM capacity scales with the ratio of two time constants, one characterizing the synaptic depression and the other one synaptic current decay time, while also increasing with facilitating time constant in a weaker way via logarithmic term, and (2) the capacity is controlled by the background excitation that should be above the critical level below which no items can be maintained in WM. The second conclusion was already illustrated in the simulations presented above. To test the validity of the first



**Figure 2. WM Capacity: Calculations and Numerical Results**

(A) Activity of two consecutively activated excitatory clusters,  $\mu$  and  $\mu + 1$ , and the inhibitory cluster. The cluster  $\mu$  elicits a PS at time  $t_0$  (magenta trace), which triggers the response of the inhibitory neural pool (black trace), the latter generating negative synaptic current at the cluster  $\mu + 1$  (green trace). The cluster  $\mu + 1$  then recovers from inhibition until its current reaches the threshold point where next PS is emitted.

(B) The function  $F(h)$  that determines the simplified one-dimensional dynamics of synaptic current according to Equation 7. Synaptic current recovers from initial hyperpolarized level  $h_0$  to the slowest point of the flow at  $h_{\min}$ , after which the flow speed is increasing and a PS is generated. The speed of the synaptic current flow is illustrated with the red arrows.

(C) The WM capacity as a function of  $\tau$ , obtained with numerical simulations of the model.  $\tau_d$  and  $\tau_f$  as in Figure 1.

(D) Same as (C); analytically estimated with Equations 5, 6, and 9. In Equation 9, we neglected the dependence of  $C$ ,  $h_0$ , and  $I_{\text{crit}}$  on synaptic parameters of the model and replaced them by the constant values  $C = 4$ ,  $h_0 = -200\text{Hz}$ , and  $I_{\text{crit}} = 2.45\text{Hz}$ . See Supplemental Information for details.  $\tau_d$  and  $\tau_f$  as in Figure 1.

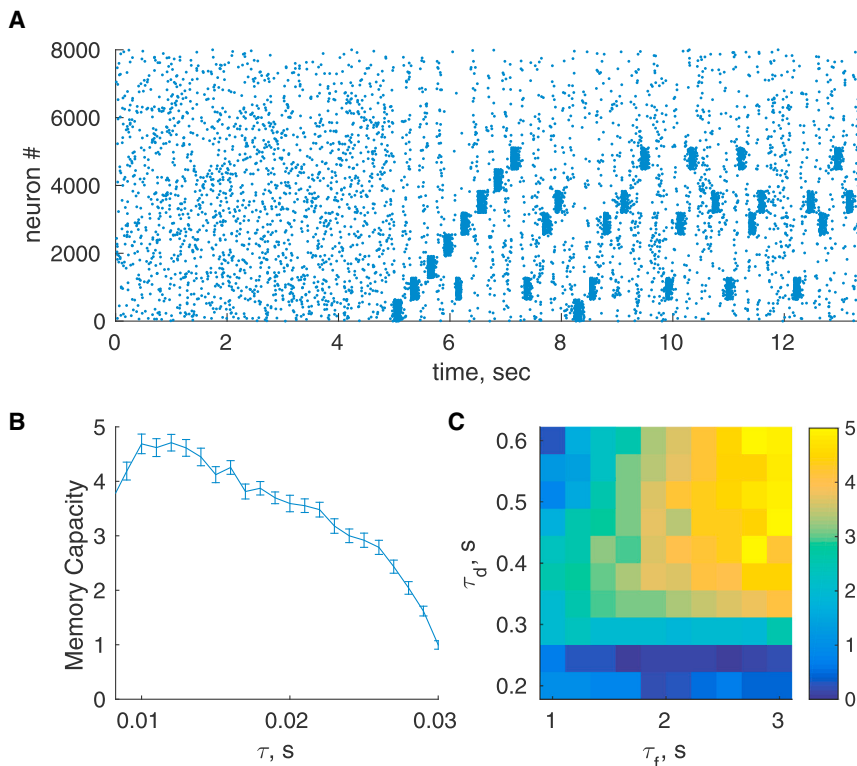
(E) The WM capacity as a function of  $\tau_f$  and  $\tau_d$  obtained with numerical simulations of the model.

(F) Same as (E); analytically estimated as in (D). The parameters  $J_{EI}$ ,  $J_{IE}$ ,  $J_{EE}$ ,  $U$ ,  $\alpha$ ,  $P$ , and  $I_b$  are the same as Figure 1D.

conclusion, we simulated the model with various choices of STP timescales  $\tau_f$ ,  $\tau_d$ , and synaptic time constant  $\tau$  (see Supplemental Information for details of simulations). The results show that our analysis captured qualitatively the WM capacity: it decreases with increasing  $\tau$ ; as a function of STP timescales, capacity is increasing with both  $\tau_d$  and  $\tau_f$  in a way consistent with analytical estimate (see Figures 2C–2F). Simulation results

are broadly compatible with analytical estimates obtained with Equation 10 (compare Figures 2C and 2E with Figures 2D and 2F), apart from a constant factor of about 2 by which analytical calculation overestimates the capacity. A closer inspection of Figure 2A reveals the origin of this discrepancy: when the synaptic current of the cluster that is about to emit the PS recovers from inhibition, the network exhibits an oscillatory activity epoch





**Figure 3. Spiking Neural Networks**

(A) An example of network simulation, including spontaneous activity and WM triggered by loading ten stimuli. Spikes of 8,000 excitatory neurons are shown as dots; neurons are arranged in order such that the first 6,400 neurons are encoding ten patterns stored in the network. Out of those, eight items are loaded sequentially into the network after 5 s of spontaneous activity. Parameters are as follows:  $\tau = 20$  ms,  $\tau_f = 3$  s,  $\tau_d = 0.6$  s, and  $U = 0.2$ .

(B) The WM capacity as a function of  $\tau$ , obtained with multiple simulations of the same network and averaging the results. Error bars: SEM.

(C) The WM capacity as a function of  $\tau_f$  and  $\tau_d$ . Parameters of the spiking network model are given in the [Supplemental Information](#).

with transient increase of inhibition that delays the onset of the PS. This effect results from the nonlinearities in network dynamics and hence is not captured by our linear approximation used in estimating the inter-PS interval  $t_s$  (see [Equation 9](#)).

The model presented above was simplified greatly to allow for analytical estimates of WM capacity. It is therefore important to address the generality of obtained results. To this end, we simulated a more realistic spiking network of 20,000 integrate and fire neurons, 16,000 excitatory, and 4,000 inhibitory ones, with noisy inputs and probabilistic synaptic structure, considered in [Mongillo et al. \(2008\)](#) (see [Supplemental Information](#) for details of the network). Each memory is represented by a group of 640 excitatory neurons, 4% of the total number. An example simulation is shown in [Figure 3A](#): the network maintains four items in WM; as opposed to the simplified model above, the corresponding populations are not activated with a constant frequency and the order of activations is variable. Overall, we found that the spiking network has a similar phase diagram with WM capacity increasing with external input up until the point where spontaneous steady state becomes unstable. By repeating simulations with different values of synaptic parameters, we observed similar trends of maximal WM capacity increasing with STP time constants  $\tau_d$  and  $\tau_f$  and decreasing with excitatory synaptic time constant  $\tau$  (see [Figures 3B and 3C](#)).

Another unrealistic feature of the model concerns the absence of overlaps between representations of different memory items. One could conjecture that overlaps in representations have similar effects to direct excitatory connections between items in our simplified model. Our analysis presented in the [Supplemental Information](#) indeed confirmed this conjecture and

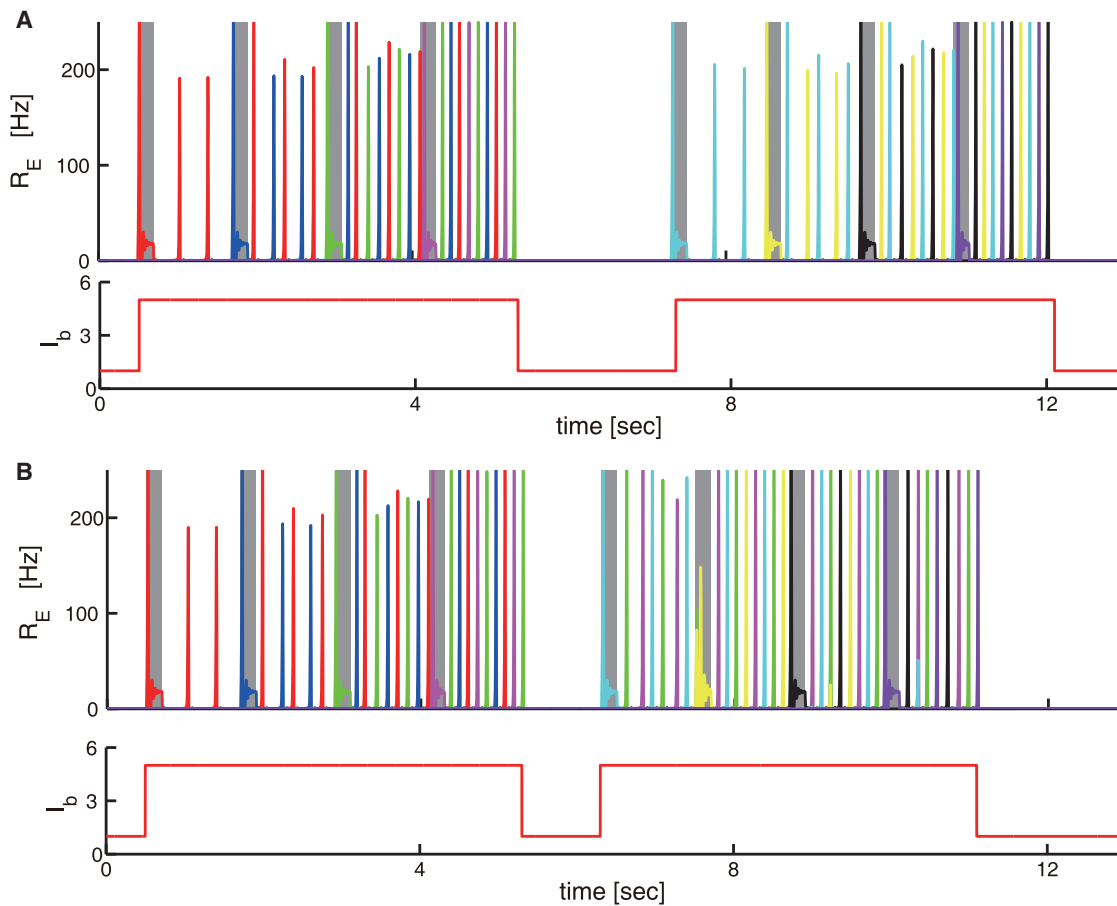
showed that the resulting model has qualitatively similar dependence of WM capacity on synaptic parameters (see [Figure S4](#)).

The dependence of WM capacity on the background excitation enables the system to be efficiently “tuned” to the desired capacity; in particular, reducing the background below the critical value removes items from WM to make room for new inputs. This tuning is crucial for many cognitive

tasks involving WM. In free recall experiments, it was shown that inserting pauses in the presentation of memory items allows human participants to encode and recall the incoming information in “chunks” of several items ([Gilbert et al., 2014](#)). We illustrate how a stream of inputs can be segmented into groups of simultaneously active items by a proper modulation of background input, if the pauses between the groups of inputs are long enough ([Figure 4A](#)). When the pauses are too short, the trace of previous items in the form of synaptic facilitation remains strong, so these items are reactivated when the background input is increased and the segmentation fails ([Figure 4B](#)).

## DISCUSSION

In this contribution, we considered the issue of short-term memory capacity in the framework of synaptic theory of WM proposed in [Mongillo et al. \(2008\)](#). We derived a simplified analytical expression for the capacity in terms of basic synaptic parameters of the network where memory items are encoded and stored in long-term memory. Surprisingly, even though the WM trace in the model is maintained by synaptic facilitation, the derived expression shows that WM capacity is chiefly increasing with the time constant of synaptic depression and only weakly increasing with the time constant of facilitation; the obtained expression also shows that capacity is inversely proportional to neuronal time constants in the network. These analytical predictions could be amenable to experimental verifications by genetic manipulations of these parameters. Our results also show that WM can be regulated by external excitation that can tune the capacity to the desired level. This regulation is crucially



**Figure 4. Segmenting a Stream of Inputs into Chunks via Modulation of Background Excitation**

(A) Two groups of inputs of four items each are loaded sequentially with a pause between them (shaded color rectangles). Upper panel: activity of each cluster is shown in a different color. Lower panel: background input. After the fourth's input is presented, all four representations from the first chunk are active until the background input is reduced. When the background input is increased again and the next four items are presented, the second chunk is active after the eighth's memory is loaded.

(B) Same as (A), but with a shorter pause between two groups of inputs. The segmentation fails because when the background input is increased after the pause, the memories from the first chunk are reactivated and mix in with the second chunk. All other parameters as in Figure 1D.

important for WM functioning, which requires not only maintaining information in a readily accessible form but also clearing it out at the appropriate time to make room for future processing. One could reasonably assume that regulation of WM reflects attentionally driven inputs to the corresponding cortical areas. Indeed, imaging studies indicate that WM tasks trigger coordination between frontal and parietal cortical regions that are considered to be involved in WM and attention, respectively, and the degree of coordination is increasing with WM load (Honey et al., 2002). We simulated the network of integrate and fire neurons introduced in Mongillo et al. (2008) and observed qualitatively similar behavior of WM capacity on synaptic parameters, confirming the generality of obtained analytical estimates.

An interesting question raised by our results is what considerations could determine a particular set of synaptic parameters in the brain that result in the experimentally observed WM capacity of approximately four for most people. More detailed analysis of the model presented in the Supplemental Information shows that

increasing the time constant of synaptic depression above a certain value brings the network to the regime where no PSs are possible and, hence, WM breaks down. The transition to this regime depends in a nontrivial manner on network parameters, but the maximal values of WM capacity tend to be substantially larger than four, both for simplified and spiking networks considered in this study. Another possibility is that WM capacity emerges from some yet unspecified functional tradeoffs. The exact nature of these tradeoffs would be difficult to determine, since networks that sustain WM could also be involved in many other cognitive processes. One potential tradeoff could be inferred from the results presented in Figure 4, where we show that in order to segment the stream of inputs into distinct chunks in WM, the pauses between the chunks should be long enough for the synaptic trace of the previous chunk to fade away, i.e., on the order of STP time constants that also determine the WM capacity. If those were made higher, segmentation of inputs would only be possible for slower presentation frequencies.



Since the WM capacity is defined by basic parameters of cortical networks, we predict that simple practice should not be enough to significantly improve the capacity beyond better tuning of attentional inputs to memory networks. In addition to having fundamental importance for understanding WM and its capacity, the analytical estimates obtained in this study might lead to new directions in clinical research of memory impairments associated with neurological disorders (Kenworthy et al., 2008; Levy and Farrow, 2001; Alloway, 2007).

## EXPERIMENTAL PROCEDURES

Please see the [Supplemental Information](#) for full experimental procedures.

## SUPPLEMENTAL INFORMATION

Supplemental Information includes Supplemental Experimental Procedures, seven figures, and one table and can be found with this article online at <http://dx.doi.org/10.1016/j.neuron.2016.12.004>.

## AUTHOR CONTRIBUTIONS

Y.M. and M.T. designed and analyzed the model, and wrote the paper. Y.M. and M.K. performed numerical simulations.

## ACKNOWLEDGMENTS

The study is supported by the EU FP7 (grant agreement 604102) and the Mortimer Zuckerman Mind Brain Behavior Institute. M.T. is also supported by Foundation Adelis. The authors thank L. Abbott and S. Fusi for helpful discussions.

Received: April 13, 2016

Revised: October 24, 2016

Accepted: November 11, 2016

Published: December 29, 2016

## REFERENCES

- Alloway, T.P. (2007). Working memory, reading, and mathematical skills in children with developmental coordination disorder. *J. Exp. Child Psychol.* 96, 20–36.
- Amit, D.J., Bernacchia, A., and Yakovlev, V. (2003). Multiple-object working memory—a model for behavioral performance. *Cereb. Cortex* 13, 435–443.
- Baddeley, A. (2003). Working memory: looking back and looking forward. *Nat. Rev. Neurosci.* 4, 829–839.
- Baddeley, A.D., and Hitch, G.J. (1974). Working memory. In *The Psychology of Learning and Motivation, Volume 8*, G.H. Bower, ed. (Academic Press), pp. 47–89.
- Colgin, L.L. (2013). Mechanisms and functions of theta rhythms. *Annu. Rev. Neurosci.* 36, 295–312.
- Compte, A., Brunel, N., Goldman-Rakic, P.S., and Wang, X.-J. (2000). Synaptic mechanisms and network dynamics underlying spatial working memory in a cortical network model. *Cereb. Cortex* 10, 910–923.
- Cowan, N. (2001). The magical number 4 in short-term memory: a reconsideration of mental storage capacity. *Behav. Brain Sci.* 24, 87–114.
- Cowan, N. (2010). The magical mystery four: how is working memory capacity limited, and why? *Curr. Dir. Psychol. Sci.* 19, 51–57.
- Edin, F., Klingberg, T., Johansson, P., McNab, F., Tegnér, J., and Compte, A. (2009). Mechanism for top-down control of working memory capacity. *Proc. Natl. Acad. Sci. USA* 106, 6802–6807.
- Fino, E., and Yuste, R. (2011). Dense inhibitory connectivity in neocortex. *Neuron* 69, 1188–1203.
- Fukuda, K., Awh, E., and Vogel, E.K. (2010). Discrete capacity limits in visual working memory. *Curr. Opin. Neurobiol.* 20, 177–182.
- Gilbert, A.C., Boucher, V.J., and Jemel, B. (2014). Perceptual chunking and its effect on memory in speech processing: ERP and behavioral evidence. *Front. Psychol.* 5, 220.
- Honey, G.D., Fu, C.H., Kim, J., Brammer, M.J., Croudace, T.J., Suckling, J., Pich, E.M., Williams, S.C., and Bullmore, E.T. (2002). Effects of verbal working memory load on corticocortical connectivity modeled by path analysis of functional magnetic resonance imaging data. *Neuroimage* 17, 573–582.
- Horn, D., and Opher, I. (1996). Temporal segmentation in a neural dynamic system. *Neural Comput.* 8, 373–389.
- Kenworthy, L., Yerys, B.E., Anthony, L.G., and Wallace, G.L. (2008). Understanding executive control in autism spectrum disorders in the lab and in the real world. *Neuropsychol. Rev.* 18, 320–338.
- Levy, F., and Farrow, M. (2001). Working memory in ADHD: prefrontal/parietal connections. *Curr. Drug Targets* 2, 347–352.
- Lisman, J.E., and Idiart, M.A. (1995). Storage of  $7 \pm 2$  short-term memories in oscillatory subcycles. *Science* 267, 1512–1515.
- Luck, S.J., and Vogel, E.K. (1997). The capacity of visual working memory for features and conjunctions. *Nature* 390, 279–281.
- Lundqvist, M., Herman, P., and Lansner, A. (2011). Theta and gamma power increases and alpha/beta power decreases with memory load in an attractor network model. *J. Cogn. Neurosci.* 23, 3008–3020.
- Lundqvist, M., Rose, J., Herman, P., Brincat, S.L., Buschman, T.J., and Miller, E.K. (2016). Gamma and beta bursts underlie working memory. *Neuron* 90, 152–164.
- Markram, H., Wang, Y., and Tsodyks, M. (1998). Differential signaling via the same axon of neocortical pyramidal neurons. *Proc. Natl. Acad. Sci. USA* 95, 5323–5328.
- Miller, G.A. (1956). The magical number seven plus or minus two: some limits on our capacity for processing information. *Psychol. Rev.* 63, 81–97.
- Miller, G.A., Galanter, E., and Pribram, K.H. (1960). *Plans and the Structure of Behavior* (Holt, Rinehart and Winston, Inc.).
- Mongillo, G., Barak, O., and Tsodyks, M. (2008). Synaptic theory of working memory. *Science* 319, 1543–1546.
- Oberauer, K. (2002). Access to information in working memory: exploring the focus of attention. *J. Exp. Psychol. Learn. Mem. Cogn.* 28, 411–421.
- Raffone, A., and Wolters, G. (2001). A cortical mechanism for binding in visual working memory. *J. Cogn. Neurosci.* 13, 766–785.
- Rolls, E.T., Dempere-Marco, L., and Deco, G. (2013). Holding multiple items in short term memory: a neural mechanism. *PLoS ONE* 8, e61078.
- Siegel, M., Warden, M.R., and Miller, E.K. (2009). Phase-dependent neuronal coding of objects in short-term memory. *Proc. Natl. Acad. Sci. USA* 106, 21341–21346.
- Tsodyks, M. (2004). Activity dependent transmission in neocortical synapses. In *Methods and Models in Neurophysics: Lecture Notes of the Les Houches Summer School 2003*, C. Chow, B. Gutkin, D. Hansel, C. Meunier, and J. Dalibard, eds. (Elsevier), pp. 245–265.
- Tsodyks, M., Pawelzik, K., and Markram, H. (1998). Neural networks with dynamic synapses. *Neural Comput.* 10, 821–835.
- Tsodyks, M., Uziel, A., and Markram, H. (2000). Synchrony generation in recurrent networks with frequency-dependent synapses. *J. Neurosci.* 20, RC50.
- Wang, Y., Markram, H., Goodman, P.H., Berger, T.K., Ma, J., and Goldman-Rakic, P.S. (2006). Heterogeneity in the pyramidal network of the medial prefrontal cortex. *Nat. Neurosci.* 9, 534–542.
- Wei, Z., Wang, X.J., and Wang, D.H. (2012). From distributed resources to limited slots in multiple-item working memory: a spiking network model with normalization. *J. Neurosci.* 32, 11228–11240.
- Zhang, M., Wang, X., and Goldberg, M.E. (2014). A spatially nonselective baseline signal in parietal cortex reflects the probability of a monkey's success on the current trial. *Proc. Natl. Acad. Sci. USA* 111, 8967–8972.