# Simulations of working memory spiking networks driven by short-term plasticity

Gianmarco Tiddia[1,2], Bruno Golosio[1,2]*, Viviana Fanti[1,2] and Pier Stanislao Paolucci[3]

[1]Department of Physics, University of Cagliari, Monserrato, Italy, [2]Istituto Nazionale di Fisica Nucleare (INFN), Sezione di Cagliari, Monserrato, Italy, [3]Istituto Nazionale di Fisica Nucleare (INFN), Sezione di Roma, Rome, Italy

Working Memory (WM) is a cognitive mechanism that enables temporary holding and manipulation of information in the human brain. This mechanism is mainly characterized by a neuronal activity during which neuron populations are able to maintain an enhanced spiking activity after being triggered by a short external cue. In this study, we implement, using the NEST simulator, a spiking neural network model in which the WM activity is sustained by a mechanism of short-term synaptic facilitation related to presynaptic calcium kinetics. The model, which is characterized by leaky integrate-and-fire neurons with exponential postsynaptic currents, is able to autonomously show an activity regime in which the memory information can be stored in a synaptic form as a result of synaptic facilitation, with spiking activity functional to facilitation maintenance. The network is able to simultaneously keep multiple memories by showing an alternated synchronous activity which preserves the synaptic facilitation within the neuron populations holding memory information. The results shown in this study confirm that a WM mechanism can be sustained by synaptic facilitation.

KEYWORDS

computational neuroscience, spiking neural networks, NEST simulation, working memory, short-term plasticity (STP)

## 1. Introduction

Working memory (WM) is a cognitive process that is able to hold and manipulate information for a short time. It is involved in a vast number of cognitive tasks (Miller et al., 1960; Baddeley and Hitch, 1974; Cowan, 1998; Golosio et al., 2015) which span from speech to visual and spatial processing. However, the WM capacity, i.e., the ability to hold multiple memories at the same time, is limited to a few items depending on the type of information (Miller, 1956; Cowan, 2001, 2010). Different from long-term memory, WM is a transient phenomenon, and it is also believed that it does not entail structural changes to the network.

A classic procedure for studying WM relies on the so-called delay response tasks. In such a framework, a stimulus is presented for a short time and the related execution of the task can take place only after a delay period. During the delay period, it is experimentally observed, especially in the prefrontal cortex (PFC), a neuronal selective

spiking activity able to maintain the information previously presented by the stimulus (Funahashi et al., 1989; Goldman-Rakic, 1995; D'Esposito and Postle, 2015). When this activity is somehow suspended (e.g., because of a noise stimulus during the delay period or a too long delay), the task is not correctly executed.

The first computational models assumed that this peculiar activity could be entirely maintained with prior long-term synaptic modifications so that when a stimulus was given to the network, the population encoding for the presented stimulus exhibited persistent spiking activity (Hebb, 1949; Hopfield, 1982; Brunel, 2000). Thus, according to these models, the information was only stored in the spiking activity. However, experimental evidence shows that memory can be also maintained when the enhanced activity is interrupted, suggesting that information is not only stored in the population's spiking activity (Stokes, 2015) but also that WM processes can exhibit discrete periodic bursts instead of a persistent activity (Honkanen et al., 2014; Lundqvist et al., 2016).

In this framework, many studies were conducted to enlighten the role of synaptic plasticity in WM (Barak and Tsodyks, 2014), and some of the proposed models rely on short-term synaptic plasticity, especially on short-term facilitation (Barak and Tsodyks, 2007; Mongillo et al., 2008; Hansel and Mato, 2013; Rolls et al., 2013). Indeed, it has been observed that the PFC shows marked short-term facilitation (Wang et al., 2006), suggesting that this form of plasticity can have a significant link with WM tasks. The work of Rolls et al. (2013) shows that employing synaptic facilitation enables a spiking network to maintain a relevant number of memories at the same time, whereas the same network lacking this kind of plasticity can maintain far fewer memories. Moreover, in Hansel and Mato (2013), it is described that the non-linearity of short-term facilitation is essential for displaying a reasonable persistent activity able to retain memory during a delay period. One of the models that posit a dominant role for synaptic facilitation in WM is Mongillo et al. (2008) model, which shows that a spiking network with synaptic facilitation is able to exhibit a bi-stable regime in which it can autonomously retain memories with periodic spiking activity without a significant firing rate increase. Thus, according to this model, memories are stored in a synaptic fashion, with spiking activity functional for synaptic facilitation upkeep. The model is further developed by Mi et al. (2017) to study how WM capacity can be modulated by short-term synaptic plasticity and the network's external excitation.

In addition, Mongillo et al. (2008) presented a simple mean-field model describing the firing rate behavior of an excitatory population modulated by short-term plasticity. This model has also been explored in Cortes et al. (2013), in which the short-term synaptic plasticity can lead to irregular and chaotic dynamics, facilitating transitions between network states and thus being one possible mechanism responsible for complex dynamics in cortical activity. Furthermore, Taher

et al. (2020) developed a neural mass model with short-term synaptic plasticity based upon the dynamics of a network of quadratic integrate and fire (QIF) neurons. Interestingly, this study was able to qualitatively reproduce the results of Mongillo et al. (2008) using facilitated synapses. Also, the maintenance of multiple memories was explored and presented an analytic expression for the WM capacity based on the work of Mi et al. (2017), in agreement with the value observed in the simulations of the network of QIF neurons.

Recently, Fiebig and Lansner (2016) (refer to also Fiebig et al., 2020) proposed a spiking network model based on a fast expression of Hebbian plasticity, in which memory is retained by oscillatory bursts. The authors of the above-mentioned study proposed a synaptic plasticity model based on a Hebbian learning rule supplemented by a short-term plasticity mechanism. This kind of implementation can enable a network to learn new memory representations, whereas using non-Hebbian plasticity needs prior long-term network training.

In this work, we implement the spiking network model described in Mongillo et al. (2008) using the NEST simulator. We show that the network exhibits totally comparable results with respect to the original article, underlining the role of short-term synaptic plasticity in WM tasks. The memory specific response of the network can be regulated by modulating the spontaneous activity. Moreover, the network is capable of maintaining multiple items at the same time and the number of items that can be maintained can be tuned by changing the short-term plasticity parameters.

## 2. Materials and methods

### 2.1. Short-term synaptic plasticity

In this section, short-term plasticity and its phenomenological description are introduced. For further details, please refer to Markram et al. (1998), Tsodyks et al. (1998, 2000), Barak and Tsodyks (2007), and Barri and Mongillo (2022).

Short-term plasticity (STP) is a mechanism in which the synaptic efficacy temporarily changes with a timescale on the order of hundreds or thousands of milliseconds. This phenomenon is regulated by the amount of synaptic resources (i.e., the neurotransmitters) available in the synapse at the moment of spike emission and by the calcium levels in the presynaptic terminal.

Indeed, the amount of neurotransmitters a synapse can contain is limited, and the emission of a spike diminishes the number of neurotransmitters available in the presynaptic terminal for further stimulation. Without synaptic activity, the amount of available neurotransmitters in the presynaptic terminal returns to its baseline level. This mechanism is called short-term depression (STD). Moreover, the spike arrival at

the presynaptic terminal elicits an influx of calcium ions that is responsible for the release of the vesicles in which neurotransmitters are stored. Higher calcium concentration in the terminal leads to a higher fraction of neurotransmitters released. This mechanism is called short-term facilitation (STF). The neurotransmitter release is then followed by a mechanism of calcium removal from the presynaptic terminal to restore its baseline concentration.

The coupling of these two phenomena leads to a temporary modulation of the synaptic efficacy (i.e., short-term plasticity), which can show STD-dominated or STF-dominated behaviors. The former can be observed when the mechanism of neurotransmitter restoration is slower with respect to the mechanism of residual calcium removal after spike emission and vice versa. To give a phenomenological description of STP, we can define $\tau_d$ as the time constant of the process of neurotransmitter restoration and $\tau_f$ the time constant for the calcium removal mechanism. Thus, we can observe STD-dominated dynamics when $\tau_d > \tau_f$ and STF-dominated dynamics when $\tau_d < \tau_f$.

The synaptic efficacy modulation led by STP can be described by the following phenomenological model: let $x$ be the normalized amount of available resources into the presynaptic terminal and let $u$ be the fraction of resources used in a spike emission. The spike arrival to the synaptic terminal rises the variable $u$ by a quantity $U(1 - u)$ (so that $u$ remains normalized), and the amount of resources released is equal to $ux$. Considering a synapse connecting the presynaptic neuron $i$ and the postsynaptic neuron $j$, this dynamics can be described by the following equations (Mongillo et al., 2008):

$$\begin{aligned} \frac{du_{i,j}}{dt} &= -\frac{u_{i,j} - U}{\tau_f} + U(1 - u_{i,j}) \sum_s \delta(t - t_s^{(i)}) \\ \frac{dx_{i,j}}{dt} &= \frac{1 - x_{i,j}}{\tau_d} - u_{i,j}x_{i,j} \sum_s \delta(t - t_s^{(i)}) \end{aligned} \tag{1}$$

where $\delta(\cdot)$ is the Dirac delta function and the sum is over the spike times $t_s^{(i)}$ of the presynaptic neuron $i$. The synaptic modulation takes place during the spike emission, so that

$$J_{i,j}(t) = J_{i,j}^{(abs)} u_{i,j}(t - \hat{\delta}_{i,j}) x_{i,j}(t - \hat{\delta}_{i,j}) \tag{2}$$

where $J_{i,j}^{(abs)}$ is the absolute synaptic efficacy for the synapse connecting neurons $i$ to neuron $j$ and $\hat{\delta}_{i,j}$ is the synaptic delay. Thus, when a spike is fired, the synaptic efficacy is described by the product $Jux$.

## 2.2. Spiking network model

This section describes the spiking network model implemented in this work, following the Supplementary material of Mongillo et al. (2008).

The network is composed of $N_E$ excitatory and $N_I$ inhibitory leaky integrate-and-fire (LIF) neurons with exponential

postsynaptic currents. The sub-threshold dynamics of the LIF neuron model is described by the differential equation

$$\tau_m \frac{dV_j}{dt} = -V_j + R_m(I_j^{exc} + I_j^{inh} + I_{ext,j}) \tag{3}$$

where $\tau_m$ is the membrane time constant, $V_j$ is the neuron's membrane potential, $R_m$ is the membrane resistance, $I_j^{exc}$ and $I_j^{inh}$ represent the excitatory and inhibitory synaptic currents received as input from the connections within the other neurons of the network and $I_{ext,j}$ represents the external input to the network.

The network external input is modeled with Gaussian white noise currents defined by the following

$$I_{ext,j}(t - \hat{\delta}_j) = \mu_{ext} + \sigma_{ext}G_k \text{ for } k\Delta t_{ng} \leq (t - \hat{\delta}_j) \leq (k+1)\Delta t_{ng} \tag{4}$$

In particular, the noise is approximated by a piecewise constant current with mean $\mu_{ext}$ and standard deviation $\sigma_{ext}$, with constant current during time intervals of length $\Delta t_{ng} = 1$ ms. Denoting the index of the time interval with $k$, for each interval, the current is given by $\mu_{ext} + \sigma_{ext}G_k$, with $G_k$ a random number extracted from a standard Gaussian distribution. The term $\hat{\delta}_j$ indicates the delays.

The synaptic current shown in Equation (3) is the sum of the contributions given by the connections with the neurons of the network, and it is characterized by excitatory and inhibitory contributions defined as $I_j^{exc}(t)$ and $I_j^{inh}(t)$, respectively. Thus, the synaptic input for a neuron $j$ of the network, with exponential postsynaptic currents, is given by the following equations for excitatory and inhibitory currents, respectively:

$$\begin{aligned} \tau_{exc} \frac{dI_j^{exc}}{dt} &= -I_j^{exc} + \sum_i \alpha J_{i,j}(t) \sum_s \delta(t - t_s^{(i)} - \hat{\delta}_{i,j}) \\ \tau_{inh} \frac{dI_j^{inh}}{dt} &= -I_j^{inh} + \sum_i \alpha J_{i,j} \sum_s \delta(t - t_s^{(i)} - \hat{\delta}_{i,j}) \end{aligned} \tag{5}$$

where $i$ is the index of the presynaptic neurons targeting the neuron $j$. $\tau_{exc}$ and $\tau_{inh}$ represent the time constant of the excitatory and inhibitory synaptic currents, respectively. In this model, $\tau_{exc} = \tau_{inh} = 2$ ms. $\hat{\delta}_{i,j}$ is the synaptic delay for the synapse connecting neurons $i$ and $j$. All the delays are uniformly distributed between 0.1 and 1.0 ms. The time dependence of the synaptic efficacy $J_{i,j}$ is only due to short-term plasticity modulation, and it is described in Equation (2), whereas synapses not modulated by the STP dynamics have fixed values of $J_{i,j}$. Since in this model, only the connections between excitatory neurons employ short-term plasticity, the connections with inhibitory neurons do not show a time dependent synaptic efficacy. In addition, since the synaptic efficacies $J_{i,j}^{(abs)}$ are expressed in mV, a factor $\alpha$ is needed in order to be consistent with the units of Equation (5). This term derives the variation of current input needed to elicit a unit of variation of the

postsynaptic potential (refer to Supplementary material for its derivation).

The excitatory neurons are organized into five selective populations, each of which includes a fixed fraction of neurons, and a non-selective population that includes the rest of the excitatory neurons of the network. In the base model, the selective populations of excitatory neurons have no overlap, so a neuron cannot belong to different selective populations. However, in an extension of the model, it is possible to simulate it with overlapping selective populations. In such a framework, the neurons belonging to each selective population are randomly chosen from the whole excitatory population, enabling the possibility of having neurons belonging to more than one selective population.

Regarding network's connectivity, short-term plasticity is implemented in all the excitatory-to-excitatory connections using the same time constants in order to show synaptic facilitation. These connections are thus characterized by the STP variables $x$ and $u$ and a weight $J$, which represents the absolute synaptic efficacy. The weights of the connections within excitatory neurons belonging to the same selective population assume a potentiated value $J_p$, emulating the result of prior long-term Hebbian learning. On the other hand, connections between excitatory neurons belonging to different selective populations, or linking a selective population with the non-selective one, are set to a baseline value $J_b$. The rest of the excitatory-to-excitatory connections have the baseline synaptic efficacy except for the 10% of them that show the potentiated value. While the excitatory-to-excitatory connections show STP dynamics, the other connections are static connections with hard-coded synaptic weights. The overall connectivity is structured so that each neuron of the network receives a fixed amount of connections from the network's populations, both excitatory and inhibitory, with non-specific inhibitory connectivity. The possibility of having more than one synapse with the same two neurons is also enabled. A simplified scheme of the spiking network architecture is depicted in Figure 1.

Such a network is able to store memories by exploiting the STP mechanism of the excitatory-to-excitatory connections. In fact, when a signal targets a selective population increasing its spiking activity, the synapses connecting neurons of the targeted population remain facilitated for a time in the order of $\tau_f$. The connections between neurons inside the pre-stimulated selective population are potentiated because $J^{(abs)}$ is relatively large due to prior long-term Hebbian learning (having $J^{(abs)} = J_p$) and also because of the STP modulation driven by $u$ and $x$. In particular, the variable $u$ shows a slow decay to its baseline value, whereas $x$ grows rapidly toward its asymptotic value. On the other hand, connections between neurons belonging to other selective populations are relatively weaker because they lack short-term potentiation, while the connections between neurons belonging to different selective populations are weaker because they have not been previously potentiated by long-term Hebbian learning (with $J^{(abs)} = J_b$ in this case). A similar effect could be driven by a random fluctuation of neuron activity, inducing an autonomous winner-take-all (WTA) mechanism, according to which the selective excitatory population with the highest firing rate stimulates the inhibitory population eliciting a suppression of the spiking activity of the other excitatory populations (Coultrip et al., 1992) due to the global inhibition. Indeed, the global inhibition is granted by the non-specific inhibitory connectivity, in agreement with experimental observations (Fino and Yuste, 2011). This mechanism decreases the amount of available resources $x$ and increases the value of $u$ across the pre-stimulated selective population. When $x$ returns close to its baseline value, since the value of $u$ is still relatively high, the connection strength becomes large enough to trigger again the WTA mechanism. This process can be reactivated periodically, and the period of reactivation is related to the dynamics of $x$ and in particular to the time constant of the synaptic depression $\tau_d$.
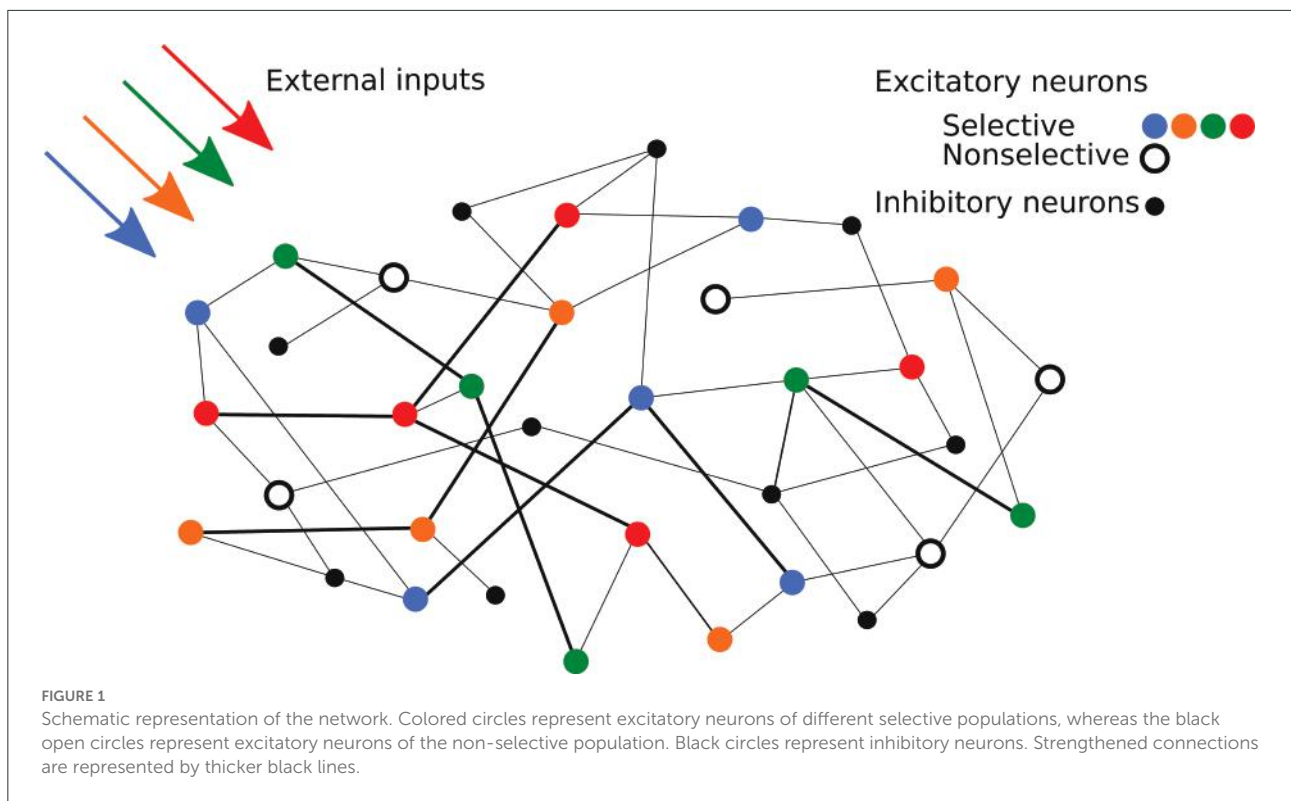
All the parameters used for the spiking network simulations, together with an in-depth description of the network connectivity, are reported in the Supplementary material.

## 3. Results

In this section, we present the results of the spiking network simulations performed using the NEST simulator (version 3.1) (Deepu et al., 2021).

The network is composed of 8,000 excitatory and 2,000 inhibitory LIF neurons with exponential postsynaptic currents, whose dynamics are described by Equations 3, 4, and 5 [refer to also Equations 1, 2, 4, and 5 in Burkitt (2006) and Equation 3 in Hanuschkin et al. (2010)]. The neuron model differs from the one employed in the original work, as in Mongillo et al. (2008) a LIF neuron model with the instantaneous rise and decay times for postsynaptic currents is employed. As discussed in Section 2.2, the excitatory population is further divided into five selective populations of 800 neurons each and a non-selective population that includes the rest of the excitatory neurons. All excitatory-to-excitatory connections follow an STP dynamics whereas the rest of the connections have fixed synaptic efficacies.

The simulations are performed using a time step of 0.05 ms, with the system of Equations (3) and (5) integrated following the exact integration scheme of Rotter and Diesmann (1999) and assuming that the external current $I_{ext,j}$ is a piecewise constant over time intervals of width $\Delta t_{ng}$. This is an additional difference with respect to Mongillo et al. (2008), in which both Equation (3) describing the neuron sub-threshold dynamics and Equation (1) describing the STP mechanism are integrated using the Euler scheme. In the NEST implementation presented here, Equation (1) is not integrated at every time step, but the values

**FIGURE 1**
Schematic representation of the network. Colored circles represent excitatory neurons of different selective populations, whereas the black open circles represent excitatory neurons of the non-selective population. Black circles represent inhibitory neurons. Strengthened connections are represented by thicker black lines.

of the variables $x_{i,j}$ and $u_{i,j}$ are analytically obtained whenever a spike is emitted by the presynaptic neuron $i$. In particular, having two consecutive spikes emitted at times $t_s$ and $t_{s+1}$ and knowing $x(t_s)$ and $u(t_s)$, the evolution of variables is computed as follows:

$$
\begin{aligned}
x(t_{s+1}^-) &= 1 + \left(x(t_s^+) - 1\right)e^{-(t_{s+1}-t_s)/\tau_d} \\
u(t_{s+1}^-) &= U + \left(u(t_s^+) - U\right)e^{-(t_{s+1}-t_s)/\tau_f} \\
u(t_{s+1}^+) &= u(t_{s+1}^-) + U\left(1 - u(t_{s+1}^-)\right) \\
x(t_{s+1}^+) &= x(t_{s+1}^-) - u(t_{s+1}^+)x(t_{s+1}^-)
\end{aligned}
\tag{6}
$$

where $t_s$ represents the spike time, while $t_s^-$ and $t_s^+$ represent the times immediately before and immediately after the spike emission, respectively. More formally, $x(t_s^-)$ and $u(t_s^-)$ can be intended as the left-side limits:

$$
\begin{aligned}
x(t_s^-) &= \lim_{\epsilon \to 0} x(t_s - \epsilon) \quad \text{with } \epsilon \in \mathbb{R}^+ \\
u(t_s^-) &= \lim_{\epsilon \to 0} u(t_s - \epsilon) \quad \text{with } \epsilon \in \mathbb{R}^+
\end{aligned}
\tag{7}
$$

while $x(t_s^+)$ and $u(t_s^+)$ can be intended as the right-side limits:

$$
\begin{aligned}
x(t_s^+) &= \lim_{\epsilon \to 0} x(t_s + \epsilon) \quad \text{with } \epsilon \in \mathbb{R}^+ \\
u(t_s^+) &= \lim_{\epsilon \to 0} u(t_s + \epsilon) \quad \text{with } \epsilon \in \mathbb{R}^+
\end{aligned}
\tag{8}
$$

Because of the discontinuity due to the spike emission, in general, the left-side and right-side limits differ from

each other for the variables $x$ and $u$. On the other hand, the exponential functions appearing in the first two lines of Equation (6) are continuous everywhere; therefore, the left and right limits are equal to each other for these functions. Therefore, the modulation led by short-term plasticity shown in Equation (2) is given by $u(t_{s+1}^+)x(t_{s+1}^-)$, thus considering variable $x$ immediately before the spike emission and the variable $u$ updated at the time of the spike emission as described in Tsodyks et al. (1998). Only after spike emission, the variable $x$ is decreased because of neurotransmitter release. This order of update stems from the fact that the presynaptic spike triggers facilitation (i.e., the increase of the variable $u$) just before the spike emission to the postsynaptic neuron. Equation (6) is implemented in the NEST simulator with the `tsodyks3_synapse` model, a modified version of the NEST synapse model `tsodyks2_synapse` model, which describes the STP dynamics according to Equation (1) as well but modulates the synaptic efficacy using the term $u(t_{s+1}^-)x(t_{s+1}^-)$. Indeed, such a difference in the implementation can be relevant, especially with neurons having low firing rates (Gast et al., 2021), with `tsodyks3_synapse` model showing higher modulated synaptic efficacies than `tsodyks2_synapse` model (refer to Supplementary material for a comparison between the two synaptic models). In this model, the STP timescales are set so that the network shows synaptic facilitation, in fact, $\tau_d = 200$ ms and $\tau_f = 1,500$ ms in agreement with the parameters chosen in Mongillo et al. (2008).

All the simulations begin with a time period of 3,000 ms in which only the background input is injected into the whole network in order to allow the network to enter its baseline state illustrating spontaneous activity. This stimulation, as well as all the other external signals, is created using the NEST `noise_generator`, which injects a Gaussian white noise current as described in Equation (4). The background input targets both excitatory and inhibitory neurons with different mean current values. Later in this section, it will be shown how network behavior can be modulated by changing excitatory background activity.

After the network reaches its spontaneous activity, an additional current, designed as a Gaussian white noise current which sums up to the background input, is injected only into a selective population for 350 ms. As a result, an item is loaded into the model. This signal, called item loading, increases the synaptic activity of the target population and thus permits a temporary strengthening of synaptic efficacies by changing the STP variables $u$ and $x$ across the connections of neurons belonging to the target population. Thus, even after the end of the item loading signal, the loaded memory can be maintained especially because of the slow decaying dynamics of the variable $u$ due to synaptic facilitation.
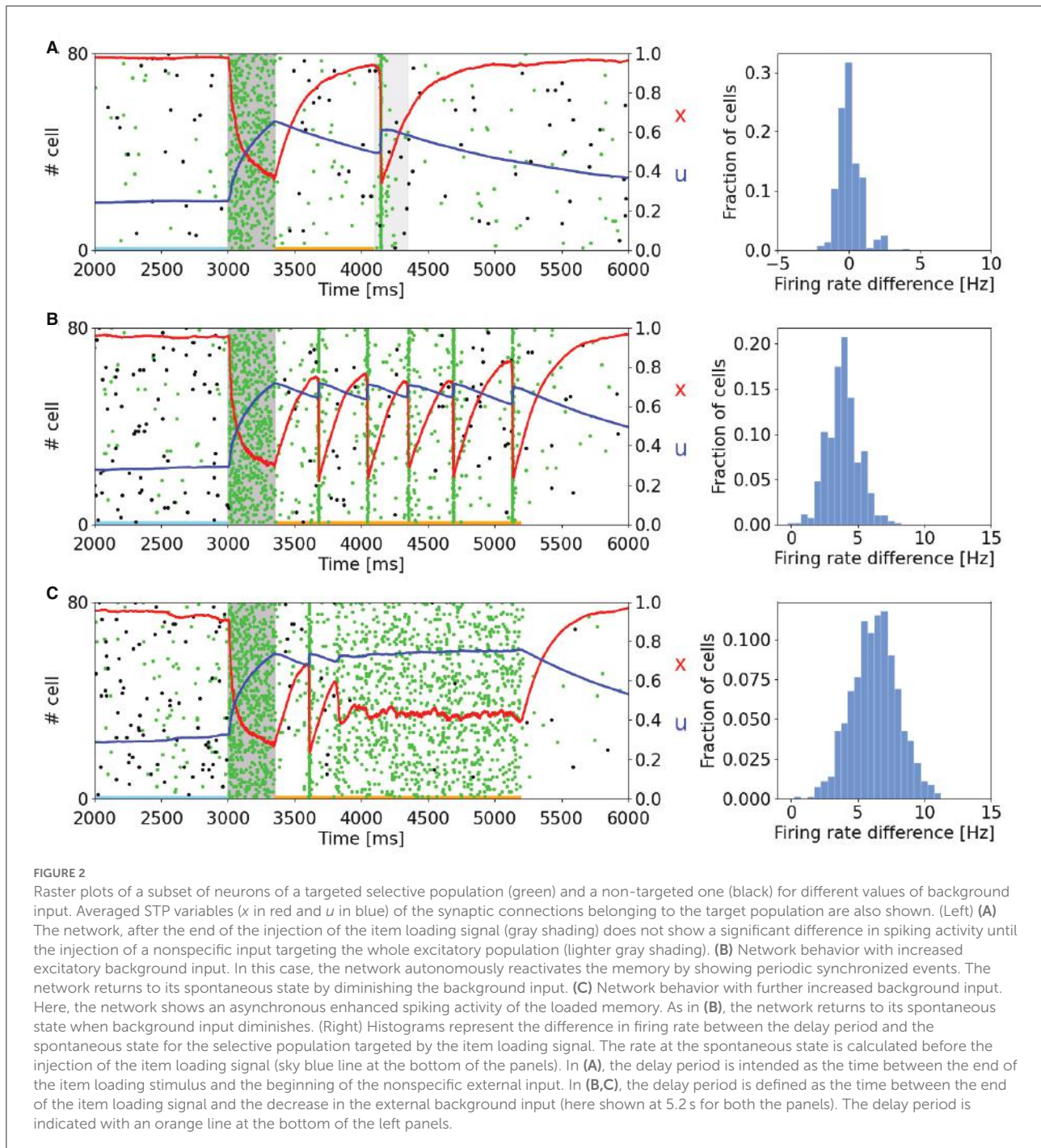
As can be seen in Figure 2, the memory specific response of the network depends on the background activity level of the excitatory neurons. This figure shows the raster plot of two selective populations, one targeted by the additional current which loads the item and a non-targeted one, together with the STP variables $x$ and $u$ averaged over the connections outbound from the neurons of the targeted selective population. In Figure 2A, to reactivate a memory, a supplemental external signal targeting the entire excitatory population is given. Although this external signal is nonspecific, only the population in which the memory was previously restored responds with the emission of a single synchronized activity, called a population spike. The network can also autonomously exhibit a memory specific spiking activity when a higher excitatory background current is injected (Figures 2B,C). In Figure 2B, the selective population which receives the item loading stimulation shows an autonomous and synchronous emission of population spikes. It should be noted that after each population spike, the STP variable $u$ increases and returns to similar values reached at the end of the item loading signal injection, interrupting the exponential decrease due to the calcium removal mechanism and thus enabling a new population spike to emerge. This behavior, together with the fast exponential growth of available resources described by the variable $x$, leads to a new stable state for the network together with the one representing spontaneous activity. To interrupt the network persistent activity, we set the excitatory background current to the value of Figure 2A. In Figure 2C the background input is further increased, and the network spontaneously shows an asynchronous higher rate activity. In this state, the memory is maintained in both spiking

and synaptic form since the STP parameters reach stable values during the high activity state followed by a population spike. As in the previously described state, the network could pass from the memory specific activity state to the spontaneous state by diminishing the background input. Indeed, without the diminishing of the background input, the network would continue to behave showing the asynchronous higher rate activity or the synchronous emission of population spikes. The values of the background current used in Figure 2 are reported in the Supplementary Table S2. Moreover, we quantitatively estimated the difference in firing rate for the targeted selective population between the delay period and the spontaneous activity state. The difference in firing rate for a neuron population is obtained by measuring the spike-count rate for each neuron of the population at two-time intervals. Naming $r_s$ the firing rate measured during the spontaneous activity state and $r_d$ the firing rate measured during the delay period, the firing rate difference for a neuron $i$ of the population is

$$\Delta r^{(i)} = r_d^{(i)} - r_s^{(i)} = \frac{N_d^{(i)}}{\Delta t_d} - \frac{N_s^{(i)}}{\Delta t_s} \qquad (9)$$
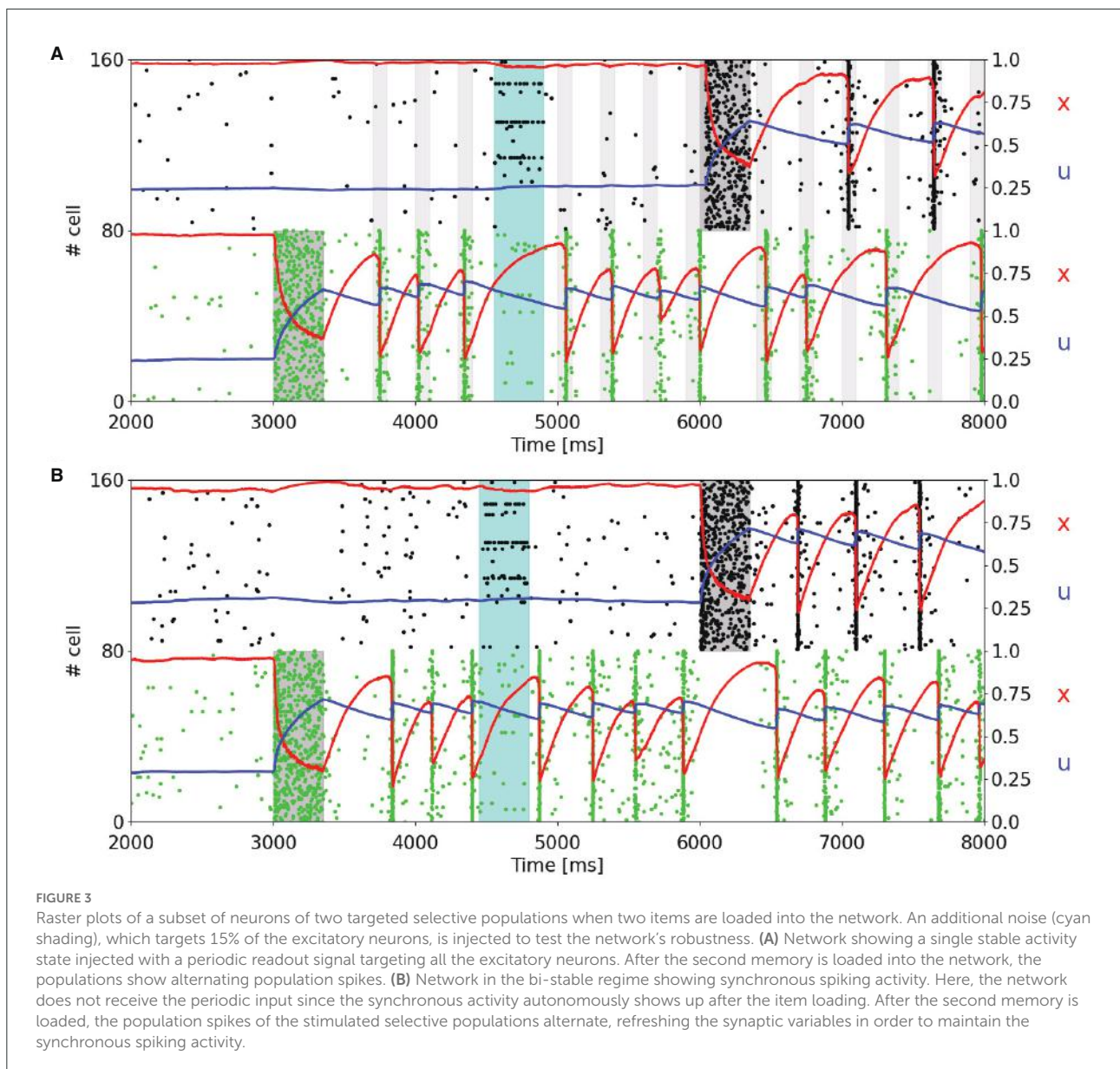
where $N^{(i)}$ is the number of spikes emitted by neuron $i$ in a certain time interval $\Delta t$. Those values are obtained for each neuron of the targeted selective population and are collected in the histograms on the right side of Figure 2. In Figure 2A, the delay period is defined as the time between the end of item loading and the beginning of the nonspecific signal, whereas in the other panels, it is identified between the end of the item loading and the decreasing of the external input (happening at 5.2 s for both panels). The time intervals related to the spontaneous activity and the delay period are indicated with horizontal lines (sky blue and orange, respectively) in the left panels of Figure 2. It is possible to notice that in Figure 2A there is no significant difference in firing rate, and a relevant part of the network shows a decrease in firing rate during the delay period. In Figures 2B,C, we observed an increase in firing rate of about 4 and 7 Hz, respectively, with an average baseline firing rate of about 0.7 Hz. These changes in firing rate are lower with respect to the ones shown in network models relying only on persistent activity to show WM behavior such as Brunel (2000) and they are in agreement with experimental measures on single-cell activity during the delay period (Shafi et al., 2007), according to which the changes in firing rate are mostly below 5 Hz.

In comparison with the work of Mongillo et al. (2008), the network simulated with NEST shows qualitatively similar results, with comparable behavior when modulating the background input targeting the excitatory neurons. However, we noticed some relevant differences with respect to the original work. For instance, on the left side of Figure 2B, it is possible to see that the time interval between adjacent population spikes is around 300 ms, whereas in Mongillo et al.

FIGURE 2
Raster plots of a subset of neurons of a targeted selective population (green) and a non-targeted one (black) for different values of background input. Averaged STP variables (x in red and u in blue) of the synaptic connections belonging to the target population are also shown. (Left) **(A)** The network, after the end of the injection of the item loading signal (gray shading) does not show a significant difference in spiking activity until the injection of a nonspecific input targeting the whole excitatory population (lighter gray shading). **(B)** Network behavior with increased excitatory background input. In this case, the network autonomously reactivates the memory by showing periodic synchronized events. The network returns to its spontaneous state by diminishing the background input. **(C)** Network behavior with further increased background input. Here, the network shows an asynchronous enhanced spiking activity of the loaded memory. As in **(B)**, the network returns to its spontaneous state when background input diminishes. (Right) Histograms represent the difference in firing rate between the delay period and the spontaneous state for the selective population targeted by the item loading signal. The rate at the spontaneous state is calculated before the injection of the item loading signal (sky blue line at the bottom of the panels). In **(A)**, the delay period is intended as the time between the end of the item loading stimulus and the beginning of the nonspecific external input. In **(B,C)**, the delay period is defined as the time between the end of the item loading signal and the decrease in the external background input (here shown at 5.2 s for both the panels). The delay period is indicated with an orange line at the bottom of the left panels.

(2008), this value is closer to 200 ms, same order of $\tau_d$. Furthermore, while the behavior of the variable $u$ is mostly comparable to the one shown in the original article, the behavior of the variable $x$ shows a considerably higher drop of the averaged variable in correspondence to a population spike. This pronounced drop in the value of $x$ is probably the reason for the difference in the time interval between the two population spikes previously mentioned.

Since one of the main features of a WM network is the holding of multiple information, we load two items into two different selective populations at different times to analyze the behavior of the STP variables of the targeted populations and the capacity of such a network of maintaining multiple items. Figure 3 shows a subset of two selective populations targeted by the item loading signal in the single stable state regime (Figure 3A) and in the regime showing synchronous and

Raster plots of a subset of neurons of two targeted selective populations when two items are loaded into the network. An additional noise (cyan shading), which targets 15% of the excitatory neurons, is injected to test the network's robustness. **(A)** Network showing a single stable activity state injected with a periodic readout signal targeting all the excitatory neurons. After the second memory is loaded into the network, the populations show alternating population spikes. **(B)** Network in the bi-stable regime showing synchronous spiking activity. Here, the network does not receive the periodic input since the synchronous activity autonomously shows up after the item loading. After the second memory is loaded, the population spikes of the stimulated selective populations alternate, refreshing the synaptic variables in order to maintain the synchronous spiking activity.

autonomous reactivation (Figure 3B), obtained using the same values of background input used in Figures 2A,B, respectively. Moreover, in both simulations noise is given to a fraction of all the excitatory neurons in order to check the robustness of the network state. The noise signal is designed as the item loading one but targets the 15% of the excitatory neurons randomly. In Figure 3A, the reactivation of the selective populations is enabled by a periodic nonspecific input (with a period of 300 ms). It can be noticed that in this framework, the two targeted selective populations do not emit the population spikes during the same periodic readout signal, but they alternate in order to reach suitable values of STP variables to enable the emission of a population spike in the following readout signal.

This peculiar behavior can also be seen when the network autonomously shows synchronous spiking activity (Figure 3B). In this case, similarly to Figure 3A, the synchronous activity of the targeted selective populations is alternated, increasing the average value of $x$ for a population when the other one is emitting the population spike. However, in Figure 3B, this mechanism is completely autonomous. In both the network states, the slow dynamics of $u$ has a key role in holding the information, in particular when another selective population shows a higher spiking activity. In addition, it can be observed that the higher spiking activity of a selective population inhibits the other populations. This is due to the network's connectivity which enables a winner-take-all mechanism, i.e., the competition

between different populations through a mechanism of global inhibition, as previously described. For this reason, it is not possible to correctly load multiple items at the same time, and it is not possible to have population spikes from different selective populations at the same time. As can be seen in Figure 3A, even if the readout signal targets all the selective populations, only the targeted selective population which has the highest STP-modulated synaptic efficacy is capable of emitting a population spike, inhibiting the excitatory neurons of the competing selecting populations.

The behavior of the network in Figure 3 is totally comparable with respect to the results shown in Mongillo et al. (2008). The main differences that emerge are related, as stated before, to the dynamics of the STP variable $x$, which shows a more pronounced drop when neurons show synchronous firing activity. We slightly increased the time interval between two consecutive stimulations in Figure 3A, from 250 to 300 ms, to make the STP variable $x$ recover enough to enable the synchronous activity response as in Figure 2A. Shorter time intervals between subsequent readout signals could result in stimulation that leads to a population spike right after the end of the stimulus.

To further test the network capacity of maintaining multiple items at the same time, we perform simulations with an additional item loading signal targeting a third selective population for the network state showing synchronous spiking activity (same value of background input as in Figure 2B). The raster plot, together with the averaged STP variables for the three targeted selective populations, is depicted in Figure 4.

As shown in Figure 4, the network is able to maintain three selective populations in the persistent activity state similar to Figure 3B such that, when all the items are loaded, population spikes alternate within the population keeping appropriate values for the STP variables. Moreover, it should be noted that each selective population in the synchronous spiking activity regime diminishes the average value for the STP variable $u$ when other items are loaded into the network. This behavior is clearly visible for the first selective population in Figure 4. Indeed, this is due to the increased distance between population spikes related to the activity of the other targeted selective populations. For an increasing number of items loaded, this can lead to a loss of synchronicity, since a persistent activity state needs to maintain a relatively high values of $u$. In this regard, we verified that an additional item loaded into the network causes the mentioned loss of synchronicity, thus such a network is able to maintain up to three items at the same time. However, an increase in the value of $\tau_f$ leads to a slower decay of the variable $u$, enabling the loading of more items into the network. For instance, with $\tau_f = 2,000$ ms, the network is able to maintain four items at the same time, and with $\tau_f = 3,000$ ms, all the five selective populations can alternate their population spikes (refer to Supplementary material). We also simulate a larger spiking network with ten selective populations in order to see whether a further increase in $\tau_f$ would enable the upkeep of more items.

We show in the Supplementary material that an opportune choice of the synaptic parameters can enable the upkeep of even seven items simultaneously (i.e., the early estimation of the WM capacity proposed by Miller, 1956). Indeed, to simulate an analogous model with ten selective populations there is a need for a larger network having a similar ratio between selective, non selective, and inhibitory populations. For this purpose, we simulated a network with 20,000 LIF neurons, with ten selective populations each composed of a similar number of neurons to that of the network described above. The parameters used to perform these simulations are presented in the Supplementary materials.

Hitherto, we presented the results of the simulations for the model with non overlapping populations, ergo an excitatory neuron can only belong to a selective population at most. To verify the network's behavior in more realistic conditions, we also performed simulations in which there is the possibility of having overlaps between the selective populations. Figure 5 shows the raster plot of a simulation with the same parameters used in Figure 3B, but with overlapping populations. Here, the population spikes are less synchronized, and not all the neurons belonging to the selective population emit a spike during the synchronous spiking activity. For this reason, the STP variable $x$ drops caused by the population spikes are less pronounced. To obtain a qualitatively similar behavior with respect to the network with non overlapping populations, the value of the potentiated synaptic efficacy $J_p$ has been slightly increased.

## 4. Discussion

In this work, we have reproduced a WM spiking network model proposed by Mongillo et al. (2008), in which the short-term synaptic plasticity has a key role in the network capability of memories upkeep. We have performed the simulations using the spiking network simulator NEST and following the network description and the parameters shown in the original work. We also modified the NEST synaptic model describing short-term plasticity to be consistent with Tsodyks et al. (1998).

Indeed, the spiking network model proposed here has some differences with respect to the original one of Mongillo et al. (2008). First, we employed a LIF neuron model with exponential postsynaptic currents, whereas the original one used a LIF model with the instantaneous rise and decay times for postsynaptic currents. Furthermore, the neuron model is integrated following the exact integration method of Rotter and Diesmann (1999), with synaptic variables for the neuron $i$ synapses updated when the neuron $i$ emits a spike, as mentioned in the Results section. The implementation of the STP dynamics follows Equations (1) and (6). In the original model, both neuron and synapse dynamics are integrated using the Euler scheme. The implementation of the STP dynamics further differs with respect to the original work. In fact, in Mongillo et al. (2008)
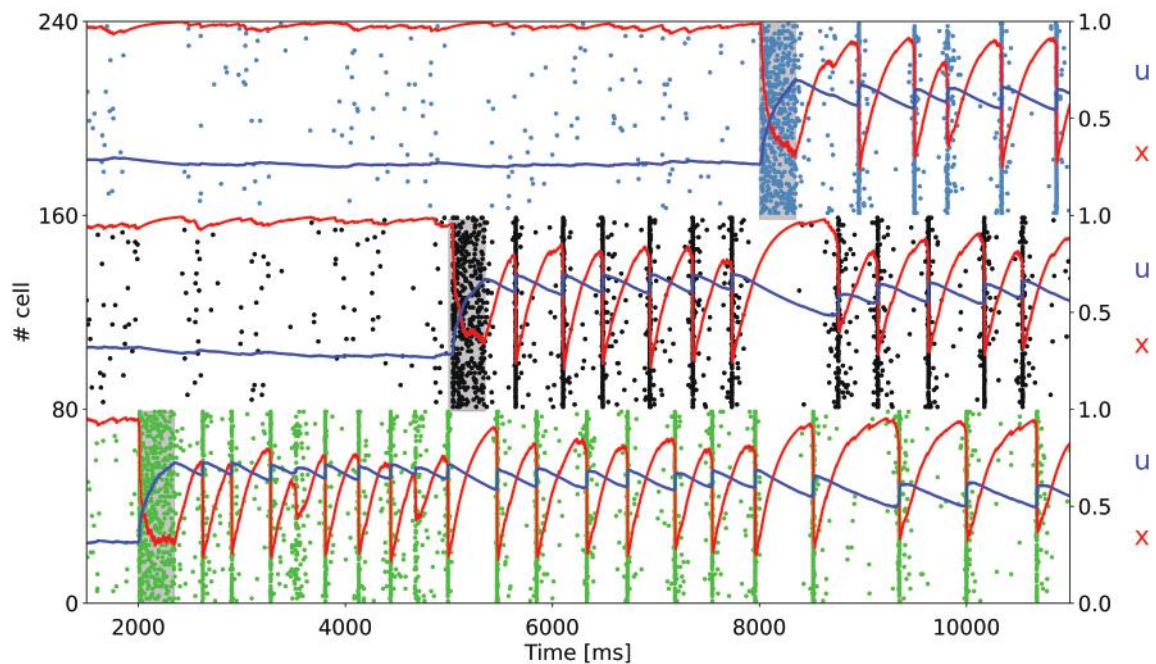
**FIGURE 4**
Raster plot of a subset of neurons of three targeted selective populations when three items are loaded into the network. The network is in the bi-stable regime showing synchronous spiking activity. Item stimuli are loaded into the network at 3,000, 6,000, and 9,000 ms.
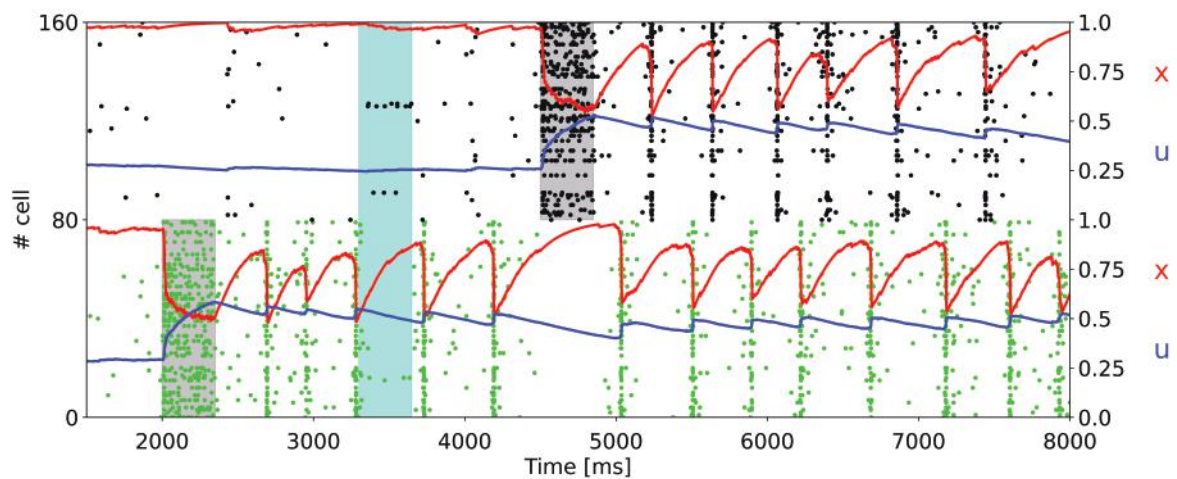


**FIGURE 5**
Raster plot for a simulation with overlapping populations. Here only a subset of the two targeted selective populations is shown. An additional noise (cyan shading) is also injected. The network is in the bi-stable regime showing synchronous spiking activity.

and Mi et al. (2017), the absolute synaptic efficacy is modulated using the values of the variables $u$ and $x$ immediately before the emission of the spikes. As described in Equation (6), the implementation used in this work considers the value of the variable $x$ immediately before the spike emission, but with the variable $u$ updated at the time of the emission of the spike, in agreement with Tsodyks et al. (1998). This change in the implementation leads to higher modulated synaptic efficacies for

the implementation employed (refer to also Gast et al., 2021 for a comparison of the two different implementations in a network of QIF neurons), and thus can be responsible for the more pronounced drop of the variable $x$ noticed in the spiking model presented in this study. Despite these differences, we were able to obtain a similar behavior with respect to the original model by slightly adjusting some parameters. However, other parameters were missing, like the integration time step. We decided to set a time step of 0.05 ms, verifying that lower or higher time steps do not entail significant changes in the network behavior (refer to Supplementary material). Moreover, the connection scheme in the NEST simulator opens to the possibility of having multiple connections with the same two neurons or also self connections. In the simulations presented here, we enabled both multiple and self connections and we verified that the dynamics of the network do not change significantly when these options are disabled (refer to Supplementary material). Furthermore, we also performed simulations with the same neuron model integrated with a different integration scheme with respect to the exact integration method of Rotter and Diesmann (1999). Specifically, we employed the stochastic Runge-Kutta method, more suitable in the presence of noise signals modeled as the background input employed in this network. We found that the results of the simulations are totally comparable with respect to the one presented in the manuscript, as shown in the Supplementary material after a description of the integration method.

Figure 2, showing the raster plot of a selective population targeted by an item loading signal and a non-targeted population, exhibits totally comparable results with respect to the original work. The response of the network can be tuned by changing the excitatory background activity, and different behaviors can thus be shown: with relatively low background activity (refer to Figure 2A), the network needs an additional excitatory signal to exhibit a memory specific response, whereas an increase in the background input can lead to a synchronous (Figure 2B) or an asynchronous (Figure 2C) higher rate persistent activity. Such responses are driven by short-term plasticity, which temporarily modulates the synaptic efficacy within the connections belonging to neurons of the targeted selective population. In particular, the slow dynamics of calcium release from the synaptic terminal grant the temporary growth of the average synaptic efficacy of the population, leading to temporary storage of the memories in a synaptic fashion. Additionally, as it is possible to see in the histograms of Figure 2, the increase in firing rate for the targeted population is relatively modest. Besides, the model can stop showing the synchronous or the asynchronous higher rate activity only by diminishing the background input current, thus was able to maintain memories for extremely long periods of time, as shown in the Supplementary material. Indeed, bifurcation analysis of a single population rate model discussed in the Supplementary material of Mongillo et al. (2008) shows that for constant inputs,

the network can show two possible behaviors: a steady state with a constant rate and a limit cycle solution corresponding to a periodic train of population spikes, which is consistent from what is observed in Figure 2B. A further study (Cortes et al., 2013) shows that, besides these stable solutions, the system can exhibit another class of states with highly irregular and chaotic-like dynamics, denoted as Shilnikov chaos.

Moreover, Figure 3 exhibits the ability of the network of maintaining multiple items at the same time and also the robustness of the network to external noise, here modeled as an item loading signal targeting a fraction of the whole excitatory neurons. It should be noted that the synchronous spiking activity of the first stimulated population, once an additional population is targeted by the item loading signal, interrupts and then alternates when the latter starts showing the synchronous persistent activity. This behavior enables the maintenance of two items at the same time. In addition, the current parameters enable the network to maintain up to three populations in the synchronous activity state (Figure 4). We verified that stimulating a fourth selective population makes the network lose synchronicity and the alternation of the so-called population spikes for the stimulated selective populations. However, we noticed that an increase in $\tau_f$ and, more generally, a change in the synaptic parameters enable the upkeep of a higher number of items at the same time (refer to Supplementary material). Indeed, as observed in Mi et al. (2017), in which a spiking network developed using the same framework as Mongillo et al. (2008) is simulated, the number of items that can be stored into the WM network (i.e., the WM capacity) can be modulated by a different choice of the synaptic parameters and of the background input. Additionally, the network has been simulated with partially overlapped selective populations to show the network's behavior in a more realistic condition. Figure 5 shows a totally comparable behavior with respect to an analog simulation with non overlapped populations, except for the fact that not every neuron of the selective population emits a spike during a population spike. Also, the population spike shows less synchronized spiking activity.

Regarding WM capacity, Mi et al. (2017) provides an analytical expression for estimating the maximum number of items that can be maintained in WM (see also Taher et al., 2020 for a similar derivation). This number is determined by the ratio between $T_{max}$, i.e., the maximal period of the limit cycle of the network and $t_s$, i.e., the time interval between two successive population spikes. Indeed, the maximal period of the limit cycle is only dependent on STP parameters and can be expressed by Mi et al. (2017)

$$T_{max} \simeq \tau_d \ln \frac{\tau_f/\tau_d}{1 - U} \qquad (10)$$

Using the parameters employed to produce Figure 4, in which up to three items can be stored in WM, $T_{max} \simeq 445$ ms, whereas the time separation between the population spike,

once the three items have been loaded into the network, is approximately $t_s \simeq 160$ ms, with the ratio between these times being

$$N_c \approx T_{max}/t_s \simeq 2.8$$

not far from the number of items stored at the same time, confirming the generality of the analytical estimation proposed in Mi et al. (2017). We also performed a similar calculation for a network able to store up to seven memories in the Supplementary material, with comparable results.

Indeed, the model presented in this work is consistent with several experimental observations. For instance, in Wolff et al. (2015, 2017) showed that, during the delay period, the information held in memory can be reactivated by a non-specific stimulus (as in Figure 2A). This result is also shown by Rose et al. (2016), in which transcranial magnetic stimulation produced a brief reactivation of the held item. Moreover, the silent dynamics can lead to interference between information from different trials (Kilpatrick, 2018), and the relation between STP dynamics and the so-called serial effects in WM tasks has recently been explored in Kiyonaga et al. (2017) and Barbosa et al. (2020). Furthermore, as already mentioned in the Results section, the firing rate changes between the spontaneous state and the delay period shown in the right panels of Figure 2 are in agreement with single-cell firing rate, which is mostly below 5 Hz and only rarely can reach values greater than 10 Hz (Shafi et al., 2007). Indeed, since a higher spiking activity would be more metabolically demanding, this behavior makes the model energetically efficient highlighting the importance of activity-silent dynamics during WM tasks and also enables a multiple memory maintenance having populations emitting bursts at different times, in agreement with Lundqvist et al. (2016).

On the other hand, this model has some limitations. The main one being that it assumes a prior long-term Hebbian learning. The way items are encoded in selective populations is extremely simplified, as all the connections within the same population have equal synaptic strength. Furthermore, this value remains constant during the simulation. A more realistic model would be a combination of long-term and short-term plasticity, enabling the learning of new items.

Working memory is responsible for the brain's ability to temporarily maintain, manipulate, and integrate information from different sensory systems (auditory, visual, etc.) during the performance of a wide range of cognitive tasks, such as learning, reasoning, and language comprehension. Indeed, this mechanism is not only useful, but can have an essential role in robotics and autonomous systems in general. For instance, WM implementation could lead to autonomous systems with cognitive capabilities closer to the human ones, enabling the possibility of learning through interactions between humans, or learning from few examples

integrating information from different sensory inputs in a similar way that humans do. Recent work has shown that the use of a WM component in robotic models can be useful to emulate many human-like cognitive functions, ranging from episodic memory, imagination and planning (Balkenius et al., 2018), language development (Giorgi et al., 2021b), and language grounding into actions and perceptions in embodied cognitive architectures (Giorgi et al., 2021a).

In conclusion, in this study, we reproduced a spiking network that shows typical WM behavior driven by short-term synaptic plasticity. This mechanism leads to a robust and energetically efficient behavior since the items loaded into the network can be maintained with a relatively low change in the population firing rate. The model, developed using the well-known spiking network simulator NEST, is available on an online repository (refer to Data Availability Statement). The NEST implementation of the model can pave the way to further studies aimed at a better understanding of WM mechanisms and of the link between short-term synaptic plasticity and long-term cognitive processes such as learning (Capone et al., 2019; Golosio et al., 2021).

## Data availability statement

The original contributions presented in this study are publicly available. The Python implementation of the spiking network model can be found at https://github.com/gmtiddia/working_memory_spiking_network. The NEST version of the model used to perform the simulations, provided with the synapse model used in this work, can be found at https://github.com/gmtiddia/nest-simulator-3.1.

## Author contributions

GT wrote the manuscript. GT, BG, and PP performed the simulations. GT, BG, VF, and PP revised the manuscript. BG and PP supervised the project. All authors have read and approved the final manuscript.

## Funding

## Conflict of interest

The authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

## Publisher's note

All claims expressed in this article are solely those of the authors and do not necessarily represent those of their affiliated organizations, or those of the publisher, the editors and the reviewers. Any product that may be evaluated in this article, or claim that may be made by its manufacturer, is not guaranteed or endorsed by the publisher.

## Supplementary material

The Supplementary Material for this article can be found online at: https://www.frontiersin.org/articles/10.3389/fnint.2022.972055/full#supplementary-material

## References

Baddeley, A. D., and Hitch, G. (1974). *Working memory. volume 8 of Psychology of Learning and Motivation*. New York, NY; San Francisco, CA; London: Academic Press.

Balkenius, C., Tjøstheim, T. A., Johansson, B., and Gärdenfors, P. (2018). From focused thought to reveries: a memory system for a conscious robot. *Front. Robot. AI* 5, 29. doi: 10.3389/frobt.2018.00029

Barak, O., and Tsodyks, M. (2007). Persistent activity in neural networks with dynamic synapses. *PLoS Comput. Biol.* 3, e104. doi: 10.1371/journal.pcbi.0030104

Barak, O., and Tsodyks, M. (2014). Working models of working memory. *Curr. Opin. Neurobiol.* 25, 20–24. doi: 10.1016/j.conb.2013.10.008

Barbosa, J., Stein, H., Martinez, R. L., Galan-Gadea, A., Li, S., Dalmau, J., et al. (2020). Interplay between persistent activity and activity-silent dynamics in the prefrontal cortex underlies serial biases in working memory. *Nat. Neurosci.* 23, 1016–1024. doi: 10.1038/s41593-020-0644-4

Barri, A., and Mongillo, G. (2022). *Short-Term Synaptic Plasticity: Microscopic Modelling and (Some) Computational Implications*. Cham: Springer International Publishing.

Brunel, N. (2000). Persistent activity and the single cell frequency-current curve in a cortical network model. *Network Comput. Neural Syst.* 11, 302. doi: 10.1088/0954-898X_11_4_302

Burkitt, A. N. (2006). A review of the integrate-and-fire neuron model: I. homogeneous synaptic input. *Biol. Cybern.* 95, 1–19. doi: 10.1007/s00422-006-0068-6

Capone, C., Pastorelli, E., Golosio, B., and Paolucci, P. S. (2019). Sleep-like slow oscillations improve visual classification through synaptic homeostasis and memory association in a thalamo-cortical model. *Scient. Rep.* 9, 8990–8911. doi: 10.1038/s41598-019-45525-0

Cortes, J. M., Desroches, M., Rodrigues, S., Veltz, R., Muñoz, M. A., and Sejnowski, T. J. (2013). Short-term synaptic plasticity in the deterministic Tsodyks-Markram model leads to unpredictable network dynamics. *Proc. Natl. Acad. Sci. U.S.A.* 110, 16610–16615. doi: 10.1073/pnas.1316071110

Coultrip, R., Granger, R., and Lynch, G. (1992). A cortical model of winner-take-all competition via lateral inhibition. *Neural Netw.* 5, 47–54. doi: 10.1016/S0893-6080(05)80006-1

Cowan, N. (1998). *Attention and Memory*. Oxford: Oxford University Press.

Cowan, N. (2001). The magical number 4 in short-term memory: a reconsideration of mental storage capacity. *Behav. Brain Sci.* 24, 87–114. doi: 10.1017/S0140525X01003922

Cowan, N. (2010). The magical mystery four. *Curr. Dir. Psychol. Sci.* 19, 51–57. doi: 10.1177/0963721409359277

Deepu, R., Spreizer, S., Trensch, G., Terhorst, D., Vennemo, S. B., Mitchell, J., et al. (2021). Nest 3.1. Available online at: https://zenodo.org/record/5508805/export/hx

D'Esposito, M., and Postle, B. R. (2015). The cognitive neuroscience of working memory. *Annu. Rev. Psychol.* 66, 115–142. doi: 10.1146/annurev-psych-010814-015031

Fiebig, F., Herman, P., and Lansner, A. (2020). An indexing theory for working memory based on fast hebbian plasticity. *eNeuro* 7, ENEURO.0374-19.2020. doi: 10.1523/ENEURO.0374-19.2020

Fiebig, F., and Lansner, A. (2016). A spiking working memory model based on hebbian short-term potentiation. *J. Neurosci.* 37, 83–96. doi: 10.1523/JNEUROSCI.1989-16.2016

Fino, E., and Yuste, R. (2011). Dense inhibitory connectivity in neocortex. *Neuron* 69, 1188–1203. doi: 10.1016/j.neuron.2011.02.025

Funahashi, S., Bruce, C. J., and Goldman-Rakic, P. S. (1989). Mnemonic coding of visual space in the monkey's dorsolateral prefrontal cortex. *J. Neurophysiol.* 61, 331–349. doi: 10.1152/jn.1989.61.2.331

Gast, R., Knösche, T. R., and Schmidt, H. (2021). Mean-field approximations of networks of spiking neurons with short-term synaptic plasticity. *Phys. Rev. E* 104, 044310. doi: 10.1103/PhysRevE.104.044310

Giorgi, I., Cangelosi, A., and Masala, G. L. (2021a). Learning actions from natural language instructions using an on-world embodied cognitive architecture. *Front. Neurorobot.* 15, 626380. doi: 10.3389/fnbot.2021.626380

Giorgi, I., Golosio, B., Esposito, M., Cangelosi, A., and Masala, G. L. (2021b). Modeling multiple language learning in a developmental cognitive architecture. *IEEE Trans. Cogn. Dev. Syst.* 13, 922–933. doi: 10.1109/TCDS.2020.3033963

Goldman-Rakic, P. (1995). Cellular basis of working memory. *Neuron* 14, 477–485. doi: 10.1016/0896-6273(95)90304-6

Golosio, B., Cangelosi, A., Gamotina, O., and Masala, G. L. (2015). A cognitive neural architecture able to learn and communicate through natural language. *PLoS ONE* 10, 1–37. doi: 10.1371/journal.pone.0140866

Golosio, B., De Luca, C., Capone, C., Pastorelli, E., Stegel, G., Tiddia, G., et al. (2021). Thalamo-cortical spiking model of incremental learning combining perception, context and NREM-sleep. *PLoS Comput. Biol.* 17, 1–26. doi: 10.1371/journal.pcbi.1009045

Hansel, D., and Mato, G. (2013). Short-term plasticity explains irregular persistent activity in working memory tasks. *J. Neurosci.* 33, 133–149. doi: 10.1523/JNEUROSCI.3455-12.2013

Hanuschkin, A., Kunkel, S., Helias, M., Morrison, A., and Diesmann, M. (2010). A general and efficient method for incorporating precise spike times in globally time-driven simulations. *Front. Neuroinform.* 4, 113. doi: 10.3389/fninf.2010.00113

Hebb, D. O. (1949). *The Organization of Behavior: A Neuropsychological Theory*. New York, NY: Wiley.

Honkanen, R., Rouhinen, S., Wang, S. H., Palva, J. M., and Palva, S. (2014). Gamma oscillations underlie the maintenance of feature-specific information and the contents of visual working memory. *Cereb. Cortex* 25, 3788–3801. doi: 10.1093/cercor/bhu263

Hopfield, J. J. (1982). Neural networks and physical systems with emergent collective computational abilities. *Proc. Natl. Acad. Sci. U.S.A.* 79, 2554–2558. doi: 10.1073/pnas.79.8.2554

Kilpatrick, Z. P. (2018). Synaptic mechanisms of interference in working memory. *Sci. Rep.* 8, 7879. doi: 10.1038/s41598-018-25958-9

Kiyonaga, A., Scimeca, J. M., Bliss, D. P., and Whitney, D. (2017). Serial dependence across perception, attention, and memory. *Trends Cogn. Sci.* 21, 493–497. doi: 10.1016/j.tics.2017.04.011

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. doi: 10.1016/j.neuron.2016.02.028

Markram, H., Wang, Y., and Tsodyks, M. (1998). Differential signaling via the same axon of neocortical pyramidal neurons. *Proc. Natl. Acad. Sci. U.S.A.* 95, 5323–5328. doi: 10.1073/pnas.95.9.5323

Mi, Y., Katkov, M., and Tsodyks, M. (2017). Synaptic correlates of working memory capacity. *Neuron* 93, 323–330. doi: 10.1016/j.neuron.2016.12.004

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. doi: 10.1037/h0043158

Miller, G. A., Galanter, E., and Pribram, K. H. (1960). *Plans and the Structure of Behavior*. New York,, NY: Henry Holt and Co.

Mongillo, G., Barak, O., and Tsodyks, M. (2008). Synaptic theory of working memory. *Science* 319, 1543–1546. doi: 10.1126/science.1150769

Rolls, E. T., Dempere-Marco, L., and Deco, G. (2013). Holding multiple items in short term memory: a neural mechanism. *PLoS ONE* 8, e61078. doi: 10.1371/journal.pone.0061078

Rose, N. S., LaRocque, J. J., Riggall, A. C., Gosseries, O., Starrett, M. J., Meyering, E. E., et al. (2016). Reactivation of latent working memories with transcranial magnetic stimulation. *Science* 354, 1136–1139. doi: 10.1126/science.aah7011

Rotter, S., and Diesmann, M. (1999). Exact digital simulation of time-invariant linear systems with applications to neuronal modeling. *Biol. Cybern.* 81, 381–402. doi: 10.1007/s004220050570

Shafi, M., Zhou, Y., Quintana, J., Chow, C., Fuster, J., and Bodner, M. (2007). Variability in neuronal activity in primate cortex during working memory tasks. *Neuroscience* 146, 1082–1108. doi: 10.1016/j.neuroscience.2006.12.072

Stokes, M. G. (2015). 'activity-silent' working memory in prefrontal cortex: a dynamic coding framework. *Trends Cogn. Sci.* 19, 394–405. doi: 10.1016/j.tics.2015.05.004

Taher, H., Torcini, A., and Olmi, S. (2020). Exact neural mass model for synaptic-based working memory. *PLoS Comput. Biol.* 16, 1–42. doi: 10.1371/journal.pcbi.1008533

Tsodyks, M., Pawelzik, K., and Markram, H. (1998). Neural networks with dynamic synapses. *Neural Comput*. 10, 821–835. doi: 10.1162/089976698300017502

Tsodyks, M., Uziel, A., and Markram, H. (2000). Synchrony generation in recurrent networks with frequency-dependent synapses. *J. Neurosci*. 20, RC50-RC50. doi: 10.1523/JNEUROSCI.20-01-j0003.2000

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. doi: 10.1038/nn1670

Wolff, M., Ding, J., Myers, N., and Stokes, M. (2015). Revealing hidden states in visual working memory using electroencephalography. *Front. Syst. Neurosci.* 9, 123. doi: 10.3389/fnsys.2015.00123

Wolff, M. J., Jochim, J., Akyürek, E. G., and Stokes, M. G. (2017). Dynamic hidden states underlying working-memory-guided behavior. *Nat. Neurosci.* 20, 864–871. doi: 10.1038/nn.4546