



HHS Public Access

Author manuscript

Neuron. Author manuscript; available in PMC 2021 May 11.

Published in final edited form as:

Neuron. 2018 October 24; 100(2): 463–475. doi:10.1016/j.neuron.2018.09.023.

Working Memory 2.0

Earl K. Miller, Mikael Lundqvist, André M. Bastos

The Picower Institute for Learning and Memory, Department of Brain and Cognitive Sciences, Massachusetts Institute of Technology, Cambridge, MA 02139

Abstract

Working memory (WM) is the fundamental function by which we break free from reflexive input-output reactions to gain control over our own thoughts. It has two types of mechanisms: online maintenance of information and its volitional or executive control. Classic models proposed persistent spiking for maintenance but have not explicitly addressed executive control. We review recent theoretical and empirical studies that suggest updates/additions to the classic model. Synaptic weight changes between sparse bursts of spiking strengthen WM maintenance. Executive control acts via interplay between network oscillations in gamma (30-100 Hz) in superficial cortical layers (layers 2 & 3) and alpha/beta (10-30 Hz) in deep cortical layers (layers 5 & 6). Deep-layer alpha/beta is associated with top-down information and inhibition. It regulates the flow of bottom-up sensory information associated with superficial layer gamma. We propose that interactions between different rhythms in distinct cortical layers underlie WM maintenance and its volitional control.

In Brief:

Miller et al present a new model of working memory (WM). Synaptic weight changes between sparse spiking help strengthen WM maintenance. Interplay between alpha/beta and gamma rhythms in different cortical layers provide an infrastructure for its volitional control.

Introduction

Working memory (WM) is the “sketchpad of conscious thought”. It is the platform where we hold and manipulate thoughts and is foundational to the organization of goal-directed behavior (Chatham and Badre, 2015; Engle et al., 1999; Fuster, 1999; Goldman-Rakic, 1995; Just and Carpenter, 1992; Miller and Cohen, 2001; Vogel and Machizawa, 2004)

Starting with work by Fuster, Goldman-Rakic and others, a wealth of data has shown that neurons in higher-order cortex, including the prefrontal cortex (PFC), show “delay activity” - elevated levels of spiking during memory delays of WM tasks (Funahashi et al., 1989;

Publisher's Disclaimer: This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting proof before it is published in its final citable form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

Declaration of Interests

The authors declare no competing interests.

Fuster and Alexander, 1971a). For example, a stimulus is shown, which must be remembered over a brief (one second or more) delay. The stimulus causes increased spiking. After it is gone, neurons continue to spike, typically at a lower rate but still above baseline levels (i.e., just before the stimulus). Everything we know suggests delay activity spiking helps maintain the WM of the stimulus. We now also know that WM involves much of the cortex. It engages executive functions associated with frontal cortex as well as posterior cortical areas that help maintain specific content (Fuster, 2015; Lara and Wallis, 2015; Miller and Cohen, 2001).

But, how, exactly, does spiking do that? Under the “classic” model, delay activity reflects persistent spiking that keeps neural ensembles “online” in a continual state of activation. However, it is important to keep in mind that virtually all of the evidence for persistent spiking is based on the time-honored practice of averaging spiking over time and across trials. This was a necessity for performing statistical analyses, especially if the data was collected one neuron at a time (as it often was, prior to the advent of multi-electrode recording). But this averaging can make spiking appear persistent even though, in real time, e.g., on single trials, it is sparse (Lundqvist et al., 2016,2018, Shafi et al., 2007).

And there are issues with persistent spiking. Spikes are metabolically expensive. Memories held by persistent spiking alone can be labile because they are lost when activity is disrupted. Multiple items can be simultaneously held if each item engages non-overlapping ensembles (Almeida et al., 2015; Edin et al., 2009). But neural ensembles often have a high degree of overlap (Fusi et al., 2016; Rigotti et al., 2013; Warden and Miller, 2010). Plus, neurons optimize information when they spike sparsely and in bursts, not persistently (Lisman, 1997; Naud and Sprekeler, 2018). In other words, in the constant chatter of the brain, a brief scream is heard better than a constant whisper. Sparse spiking also allows multiple items to be multiplexed in time, preventing them from interfering with one another and simplifying the readout of WM (different ensembles shout in turn instead of mumbling on top of each other) (Bahramisharif et al., 2017; Lisman and Idiart, 1995; Lundqvist et al., 2011; Sandberg et al., 2003; Siegel et al., 2009). In fact, even sustained attention is not truly sustained. The brain samples the environment periodically (Buschman and Miller, 2010; Fiebelkorn et al., 2018; Helfrich et al., 2018; Landau and Fries, 2012; Schroeder et al., 2010; VanRullen, 2016). All this suggests that WM (and cognition in general) is more complex than a simple persistence of spiking and average spike rates.

Further, a critical aspect of WM has not enjoyed as much experimental effort as its maintenance functions. Volitional control is what makes WM special. It is the fundamental function by which our brain wrests control of behavior from the environment and turns it to our own internal goals (Goldman-Rakic, 1995). We can choose what to think about, when and if to act. Breakdown in volition is associated with psychiatric disease, like schizophrenia (Uhlhaas and Singer, 2010). Volition is, necessarily, a network phenomenon and thus not well-addressed at the single-neuron level. Network properties can be examined with multiple-electrode recordings of multiple neurons and at the level of local field potential (LFP) level, the summed activity of many neurons. During WM tasks there are LFP oscillations (i.e. synchronized activity) in the alpha/beta band (10-30Hz), gamma band (30-100Hz), and theta band (4-8 Hz) (Bahramisharif et al., 2017; Bastos et al., 2018; van

Ede et al., 2017; Honkanen et al., 2015; Howard et al., 2003; Lundqvist et al., 2016; Roux et al., 2012; Salazar et al., 2012).

The gamma band has been associated with sensory information held in WM (Bastos et al., 2018; Honkanen et al., 2015; Howard et al., 2003; Roux et al., 2012) as well as spiking carrying sensory information (Lundqvist et al., 2016; 2018). In fact, gamma power correlates with the number of objects held in WM (Howard et al., 2003; Kornblith et al., 2016; Roux et al., 2012). The alpha/beta band has been associated with top-down information (e.g., task rules) and with inhibitory functions (discussed below). It is anti-correlated with gamma. The theta band may play a role in generating irregular bursts of gamma/spiking (see below). As we will see, the interaction between these rhythms and spiking has provided insight into top-down “executive” control that gates access to WM.

To be clear, we are not suggesting that the classic model of persistent spiking is wrong. It is correct at a certain level of approximation, averaged spiking of single neurons. But a new look in more detail (e.g., on single trials) and on a network level has provided new insights. The results still point to a central role for spiking in WM. It is just that the story is more complex than previously suspected. Is not that always the case?

Persistent Problems

We recently reviewed evidence for and against persistent spiking underlying WM (Lundqvist, M., Herman, P., Miller, E.K., 2018), so we will be brief here.

The evidence associating delay interval spiking with WM maintenance is clear and unequivocal (e.g., Funahashi et al., 1989; Fuster, 1999; Fuster and Alexander, 1971; Goldman-Rakic, 1995; Miller et al., 1996; Pasternak and Greenlee, 2005; Romo et al., 1999). However, the evidence that spiking is persistent is less so. Virtually all prior studies averaged spiking over time, across trials and often across neurons recorded in different sessions. Averaging masks the details of spiking activity (Stokes and Spaak, 2016). Single-trial analyses indicate spiking is typically sparse in real time (Kucewicz et al., 2017; Lundqvist et al., 2016, 2018; Shafi et al., 2007; Stokes and Spaak, 2016).

Yes, there are examples in the literature of single neurons that seem to show persistent spiking on individual trial rasters. This suggests that at least some neurons show persistent activity. But the bulk of neurons spike sparsely in WM delays, even when spiking is averaged across trials (Cromer et al., 2010; Fujisawa et al., 2008; Hussar and Pasternak, 2012; Shafi et al., 2007). A model that only explains the properties of a small percentage of the population is not complete. In addition, those examples are almost all from single-neuron studies in which investigators (necessarily) chose to only study neurons that seemed to show a property of interest (like delay activity spiking). That, plus single-neuron examples are invariably “best of”, means that they are hardly representative of the underlying population. Further, single-neuron studies typically optimize stimulus parameters for the neuron under study thus optimizing neural activity. Under real-world conditions, however, only a tiny fraction of neurons may be operating under such ideal conditions. They are also not representative of the bulk of neurons contributing to a given function. Parsimony suggests

that the whole population of neurons contribute to behavior not just a select few operating under ideal conditions.

We are not saying that there is anything wrong with the approaches described above. They were and are essential for identifying constituent neural mechanisms (like delay activity). However, whether spiking activity is persistent vs sparse is a different level of question. It is one of how neural populations and the circuits they form contribute. This requires an approach in which neurons are sampled more randomly and in the context of the activity of other neurons so that network properties can be deduced. For this level of question, multiple-electrode studies that record activity of dozens to hundreds of neurons are better suited than single-neuron recording (Lundqvist, Herman, Miller, 2018; Miller and Wilson, 2008).

But if single neurons do not show persistence, is it possible that it can be seen on the level of populations of neurons? This possibility rests on the assumption that single neurons spike asynchronously (i.e., at different times). When combined across different neurons, spiking “fills” time, producing persistence at the population level. To test this, one needs to measure activity in local networks, not just single neurons. This can involve analysis of multiple simultaneously recorded neurons as well as local field potentials (LFPs) which provide a measure of coordinated activity of neurons within a few hundred micrometers. We recently applied this approach to examine delay activity across seven frontal cortical areas (dorsolateral PFC, ventrolateral PFC, FEF, dorsal premotor cortex, 8A, 8B, and the supplementary motor area/anterior cingulate cortex). As expanded below, this indicated that local populations of neurons are not asynchronous. Instead, there are sparse and coordinated bursts of spiking (Bastos et al., 2018; Lundqvist et al., 2016, 2018).

Of course, one could posit that if you combine enough neurons across a wide enough expanse of cortex, one can fill time with spikes. In other words, it could be that activity is persistent when combined across highly distributed networks. However, in order to evaluate extant models, the local network is critical. Much of the brain’s computations take place on a local level. The cortex is thought to be organized into local, recurrently connected clusters with shared tuning properties (Constantinidis et al., 2001; Kritzer and Goldman-Rakic, 1995) and persistent activity is typically modelled using local recurrent connectivity (Amit and Brunel, 1997; Compte et al., 2000; Durstewitz et al., 2000).

Nonetheless, we can consider global cortical activity by using techniques like EEG and fMRI. This has revealed that for extended periods of time, information held in working memory cannot be decoded from global activity. However, when the cortex is “pinged” by a task-irrelevant stimulus or by transcranial magnetic stimulation, the network “rings” back with the information held in working memory (Rose et al., 2016; Sprague et al., 2016; Stokes et al., 2013; Wolff et al., 2017). This suggests that the WM can be held in the absence of persistent spiking.

Finally, most of the evidence that WM involves simply maintaining ensembles in a persistent active state comes from relatively simple tasks in which a single item must be retained over a “blank” delay interval (with no intervening distractions, further additions to WM etc.). That favors evidence for persistence spiking by “protecting” it from events that might disrupt

it. When an interruption occurs, for example, by having the animal focus briefly on another task, delay activity can be disrupted for 100s of milliseconds without any loss of the WM items (Spaak et al., 2017; Watanabe and Funahashi, 2014). It is possible, in principle, that WM items could be switched in and out of long-term memory to bridge these gaps. But that would still require maintenance of an index to the information in long-term memory.

Another related issue is the stability of the neural code underlying WM. New sensory inputs can change the neural population code carrying WMs. This can be evaluated by testing if a decoder trained on activity at one time in the trial can decode information at other times. If not, then there has been a change in code. Even without intervening inputs, the population code changes over the memory delay (Meyers et al., 2008; Spaak et al., 2017; Stokes et al., 2013). This argues against a model of WM in which an ensemble is activated by a sensory input and then kept in that active state. Instead, WM representations are dynamic and change over time. It is possible, however, to find a linear combination of neurons that will maintain a stable code, “a stable subspace” (Murray et al., 2017). However, this has been demonstrated with “blank” delays without additional inputs or distractions. Decoders trained on time before additional inputs do not perform well following it (Parthasarathy et al., 2017).

Further, computational modeling of persistent activity using attractor dynamics suggests its limitations. Attractor dynamics are network dynamics dominated by neurons with persistent spiking. Different attractor states correspond to unique patterns of activity corresponding to different items in working memory. As long as the state is maintained, the memories are held. The problem is that attractor states are not stable when they are perturbed. They can be disrupted by a distracting input or by adding additional information to WM. For this same reason they have difficulty storing more than one WM at a time. Bump attractor models, originally proposed for visuospatial working memory, can store multiple locations if there is no overlap in their neural representations, that is, if the WMs are held by essentially different networks (Almeida et al., 2015; Edin et al., 2009). But if there is overlap, the attractor states for different WMs tend to meld into one. This is problematic for the overlapping representations seen in the PFC (the cortical area most associated with WM), at least for non-spatial information (Rigotti et al., 2013; Warden and Miller, 2010). Any universal model of WM needs to deal with overlapping representations. Otherwise, it is only a special-case model.

What is the alternative?

An alternative is a hybrid attractor-dynamic/synaptic model. Rather than persistent spiking, there are brief, sparse, bursts of spiking. WMs are held between spiking by spiking-induced changes in synaptic weights, “impressions” left in the network (Lundqvist et al., 2011a, 2012; Mongillo et al., 2008; Sandberg et al., 2003; Stokes, 2015). Wang and Goldman-Rakic and colleagues showed that spiking in the PFC can produce fast synaptic enhancement that lasts hundreds of milliseconds (Wang et al., 2006). In fact, the enhancement depends on sparse, bursty, spiking. Not only is this metabolically less expensive, it also mitigates many of the problems of persistent attractor states. Synaptic weights are less prone to interference. Because the time spent in active attractor states is kept

to a minimum, the WMs are less prone to disruption from, e.g., a new sensory input. Multiple items can be simultaneously held by multiplexing in time their brief bouts of activity. In other words, by having different ensembles active at different times, the attractor states do not interfere with each other (e.g., Siegel et al., 2009).

For example, in the Synaptic Attractor Model (SAM) ensembles have inhibitory connections with other ensembles (Lundqvist et al., 2011), a feature shared by classic models of WM (Amit and Brunel, 1997; Goldman-Rakic, 1996). Each attractor state has a limited lifetime. Thus, they are semi-stable and shut others down temporarily. The result is that each WM item is expressed in brief bouts of spiking. Based on known biophysics, the SAM predicts that in absence of bottom-up sensory inputs networks oscillate in the alpha/beta band (10-30 Hz), only occasionally spiking. When a bottom-up sensory input activates an ensemble, it temporarily oscillates in a gamma state (>30 Hz) and gives off a short burst of elevated spiking before inhibition reverts it back to the alpha/beta state and reduced spiking. The gamma bursts may be linked to underlying theta rhythms (Canolty et al., 2006; Voytek et al., 2015; Watrous et al., 2015). This could organize time-multiplexing of items (Bahramisharif et al., 2017; Fuentemilla et al., 2010; Herman et al., 2013).

The spiking induces temporary (<1 sec) changes in synaptic weights, perhaps via calcium dynamics (Lundqvist et al., 2011, 2016; Mongillo et al., 2008; Wang et al., 2006). Therefore, both spiking and short-term plasticity are thought to be mechanisms for WM storage. Brief, irregular bursts of spiking and gamma during the memory delay are needed to occasionally refresh the synaptic weight changes so that the WMs can be maintained beyond the lifetime of the synaptic weight changes.

In the model, the refresh rate is responsible for the limited capacity of WM (an average of four items)(Awh et al., 2007; Buschman et al., 2011; Cowan, 2010; Luck and Vogel, 1997). If too many items are simultaneously held, the requirement to refresh the synapses causes a build-up of interference due to competition for the limited time available for the refresh (Lundqvist et al., 2011a; Mi et al., 2017). For this reason the gamma burst rate increases with WM load (Lundqvist et al., 2016). Schizophrenic patients have lowered WM capacity, and do not demonstrate the load-dependent changes in gamma (Basar-Eroglu et al., 2007) observed in healthy subjects (Howard et al., 2003; Roux et al., 2012).

We tested this model by analyzing local field potential (LFP) and spiking from seven cortical areas (dorsolateral and ventrolateral PFC, the frontal eye fields, dorsal premotor cortex, areas 8A and 8B, and the supplementary motor area/anterior cingulate cortex) of monkeys performing several different WM tasks (Bastos et al., 2018; Lundqvist et al., 2016, 2018). These tasks involved both spatial and non-spatial WM, and different WM loads (1-3). Across all these different tasks and areas, spiking that carried information about the sensory inputs to-be-held in WM were highly associated brief bursts of narrow-band gamma oscillations, especially during the encoding of sensory information into WM (Figure 1A-C Lundqvist et al., 2016). During such gamma bursts, spiking was elevated and more informative about the contents of WM than spiking outside the bursts (Lundqvist et al., 2018). In fact, at recording sites where spiking did not carry WM information, there was little or no gamma bursting (Figure 1D). Interleaved with the gamma bursts were brief bursts

of beta and bursts that were not associated with spiking carrying WM contents. During the memory delays, the gamma bursts occurred at a lower rate but still above the baseline rate (Bastos et al., 2018; Lundqvist et al., 2016, 2018). This is consistent with the model prediction that gamma bursts/spikes are needed to refresh synaptic weight changes. The gamma bursting and associated spiking increased near the end of the delay, around the time WMs needed to be “read out” (Figure 1C).

Importantly, gamma and alpha/beta bursts were anti-correlated, like mirror images of each other (Figure 1C). This was task-related, only appearing at recording sites where spiking reflected the contents of WM (Bahramisharif et al., 2017; Lundqvist et al., 2016, 2018). The task-related anti-correlation between gamma and alpha/beta intrigued us. It occurred to us that it could be a mechanism for controlling WM storage. Gamma is associated with the spiking that holds sensory inputs in WM. If it has a push-pull relationship with alpha/beta, then gamma (and hence WM storage) can be turned on and off by lowering and raising alpha/beta, respectively. For example, turning down alpha/beta would allow gamma to be expressed and sensory inputs to be encoded in WM. Turning up alpha/beta would turn down gamma and thus clear out the WM storage.

What about alpha/beta?

The above implies that alpha/beta has an inhibitory role in WM. In visual cortex, inhibition has been linked with alpha (8-12 Hz) (Haegens et al., 2011; Jensen and Mazaheri, 2010). In prefrontal and motor cortex, inhibition is more often linked with beta (15-30 Hz). However, several studies report power modulation that spans both the alpha and beta bands (Bastos et al., 2018; Ede et al., 2011). Thus, we will group these bands together as they seem to have similar functions: providing inhibition. One exception is in parietal cortex where lower beta has been associated with WM maintenance (Kopell et al., 2011; Salazar et al., 2012).

Motor planning has similarities to WM control and may have shared evolutionary origins (Chatham and Badre, 2015). In fact, motor beta/gamma has very similar behavioral correlates as WM beta/gamma. Beta is elevated when a movement is being withheld (Donoghue et al., 1998; Feingold et al., 2015; Jha et al., 2015; Zhang et al., 2008). During movement, beta wanes and gamma waxes. Beta is then elevated after movement (Feingold et al., 2015) as if the motor plan was being cleared out. Similarly, there was increased beta in the PFC after the end of a trial, once WMs are no longer relevant (Lundqvist et al., 2018). In fact, this effect was selective to recording sites where WM information was held during the trial. Alpha/beta may also play a role in protecting WM from distractors (Bonfond and Jensen, 2012). Across virtually all of sensory cortex, gamma is associated with sensory processing and beta is anti-correlated with gamma (Bauer et al., 2006; David et al., 2015; Ede et al., 2011; Fisch et al., 2009; Fontolan et al., 2014; Zhang et al., 2008).

Inhibition is central to executive control and so is the knowledge about what needs to be controlled (Miller and Cohen, 2001). Correspondingly, beta has also been associated with the top-down information. Task rules are reflected in different patterns of beta synchrony in PFC (Buschman et al., 2012) and visual cortex (Richter et al., 2018) as if beta was helping form ensembles for the rules. Such content-specific “beta ensembles” have also been found for

other types of top-down information like learned categories (Antzoulatos and Miller, 2016; Stanley et al., 2018; Wutz et al., 2018). Thus, with the spatio-temporal pattern of beta changing with top-down information, beta's inhibitory effects can act selectively and direct the flow of sensory information.

Support for this comes from numerous studies showing that attention to sensory inputs results in increased gamma while increased alpha/beta occurs for modalities or locations that are unattended (Buffalo et al., 2011; Fries et al., 2001; Haegens et al., 2011; Jensen and Mazaheri, 2010)(van Ede et al., 2017; Leszczynski et al., 2017; Popov et al., 2017; Wolff et al., 2017). A MEG study in humans also showed that the alpha/beta in sensory cortex was anti-correlated with the locus of attention (and with gamma) and was under top-down control from frontal cortex (Popov et al., 2017). The alpha/beta was also anti-correlated with behavioral reaction time, indicating its functional relevance.

Thus, we propose dual roles for beta: inhibition and formation of ensembles for top-down information. We hypothesize that the inhibitory role for beta is a mechanism acts locally, at the level of cortical columns (Bastos et al., 2018). This local inhibition is akin to the role proposed for alpha in sensory cortex (Jensen and Mazaheri, 2010). In addition, beta rhythms have been proposed to be ideally suited for flexibly generating neural ensembles (Kopell et al., 2011; Spitzer and Haegens, 2017) with the beta rhythmic networks reaching down to the level of individual cells (Dann et al., 2016). These large-scale neural ensembles, we propose, contain the top-down knowledge required to locally deliver inhibition, and thus executive control, where and when it is needed.

In correspondence with their roles in top-down vs bottom-up functions, beta and gamma have also been associated with feedback and feedforward cortical processing. In a study using large-scale electrocorticography, Bastos et al. recorded from eight different visual areas simultaneously (Bastos et al., 2015a) as monkeys performed a visual attention task. A cortico-cortical motif emerged by analyzing all pairs of areas in relation to their anatomical pattern of feedforward/feedback connectivity. Gamma oscillations were shown to flow up the visual cortical hierarchy in a bottom-up direction. Beta oscillations flowed down the hierarchy in the top-down direction. A similar functional hierarchy was then subsequently discovered in the human visual system with MEG recordings (Michalareas et al., 2016). Causal evidence also supports these findings. Electrical micro-stimulation in V1 causes increases in gamma power in V4, an area downstream from V1 and in receipt of feedforward connections. Micro-stimulation in V4 causes increases in alpha power in V1 (van Kerkoerle et al., 2014).

Note that bottom-up gamma is not inconsistent with the idea that top-down attention often enhances gamma power and inter-area synchrony (Bastos et al., 2015a; Bosman et al., 2012; Buschman and Miller, 2007; Fries et al., 2001; Gregoriou et al., 2009). Top-down attention is often conceptualized as a "spotlight" that turns up the gain on behaviorally relevant sensory representations (Desimone and Duncan, 1995). Thus, sensory enhancement of attended items also enhances gamma. At the same time, gamma-enhancement can be controlled by beta rhythms (Lee et al., 2013). Richter and colleagues examined the trial-by-trial pattern of top-down Granger causality from parietal to visual cortex in beta with the

bottom-up Granger causality from V1 to V4 in gamma (Richter et al., 2017). The strength of top-down (parietal to visual cortex) beta synchrony predicted the strength of bottom-up (V1 to V4) gamma synchrony.

Plugging this into what we suggested above, the idea is that top-down information carried by alpha/beta rhythms could inhibit the expression of bottom-up information carried by gamma rhythms and perhaps even regulate the precise patterns of gamma synchrony that enable corticocortical communication (Fries, 2015). But how do these rhythms interact on a micro-circuit level? The answer seemed to lie in interactions between cortical layers.

Beta in deep-layer cortex interacts with gamma in superficial-layer cortex

The cerebral cortex has laminar organization. Layer 4 is the input layer (Felleman and Van Essen, 1991; Gilbert and Wiesel, 1983; Rockland and Pandya, 1979). Although the correspondence is not perfect (biology never is), the superficial layers (layers 2-3) largely contain the feedforward-projecting neurons that carry sensory information anteriorly while the deep layers (layers 5-6) contain the feedback-projecting neurons that carry the top-down information posteriorly in cortex (Markov et al., 2013a). Gamma and beta rhythms are emphasized in different cortical layers. In visual cortex, gamma is more prominent in superficial and middle layers while alpha/beta is more prominent in deep layers (Bollimunta et al., 2008; Buffalo et al., 2011; Maier et al., 2010; Smith et al., 2013; Xing et al., 2012).

To determine if this was also true in frontal cortical areas associated with WM, we recorded with “laminar” electrodes in animals performing three different WM tasks (Bastos et al., 2018). Laminar electrodes have multiple contacts along the shaft and thus allow recording from all cortical layers simultaneously.

Frontal cortex gamma power and cue-related information peaked in superficial layers while alpha/beta peaked in deep layers (Fig. 2A). WM delay interval spiking was also stronger in superficial layers (Fig 2B). This corresponds with our observations about gamma and beta rhythms. The superficial layers are the feedforward layers that carry bottom-up sensory inputs up the cortical hierarchy. Thus, it is where we would expect to find more bottom-up gamma and spiking carrying sensory information held in WM. In sensory areas, bottom-up gamma and informative spiking is typically only elevated for the duration of sensory stimulation (Buffalo et al., 2011; Fries et al., 2001). In prefrontal cortex, bursts of spiking and gamma also appear over the delay interval. This could be the result of longer time synaptic integration constants (Murray et al., 2014) brought about by superficial-layer lateral excitatory connections, (Goldman-Rakic, 1996) and more synaptic spines on pyramidal cells (Elston, 2000). Likewise, it makes sense that beta would be stronger in deep layers. Beta is associated with top-down information. The deep layers are the feedback layers that can carry top-down information from frontal cortex down the cortical hierarchy.

The pattern of influence between beta and gamma suggested a laminar-rhythmic infrastructure for control of WM storage. Granger Causality is a statistical measure of time series prediction that is indicative of functional connectivity (Bressler and Seth, 2011). It indicated that deep-layer beta oscillations regulated superficial-layer beta. The phase of

deep-layer beta oscillations, in turn, modulated the amplitude of superficial gamma (Bastos et al., 2018; Canolty et al., 2006; Colgin et al., 2009; Lakatos et al., 2005; Spaak et al., 2012). Importantly, the power of deep-layer beta was inversely correlated with superficial-layer gamma, consistent with an inhibitory role for beta (Figure 2C). Thus, coupling between deep and superficial layers may serve a control function. Increasing deep-layer alpha/beta would increase superficial layer beta. Superficial layer beta would, in turn, suppress gamma and thus the expression of bottom-up sensory information in superficial layers. This would prevent the encoding of sensory information in WM. Conversely, if deep-layer beta is reduced, there would be decreased coupling to superficial-layer beta. That would release gamma from inhibition, allowing its expression and the encoding of bottom-up information into WM. Indeed, we found that the strength of deep-layer beta coupling to superficial-layer gamma was reduced during the WM delays compared to the pre-cue baseline period (Bastos et al., 2018).

Mechanisms of gamma/beta interplay

To understand how the interplay between gamma and beta gives rise to WM control, it is important to consider their neurophysiological origins. Here, we provide a short summary (for detailed reviews see Buzsáki and Wang, 2012; Fries, 2015; Wang, 2010). Excitatory (E) and inhibitory (I) cells are densely interconnected in cortex. Fast (greater than 10 Hz) rhythms can be generated in cortex through recurrent inhibition between E cells and a variety of classes of I cells. Fast-Spiking (FS) I cells are a key player. They provide the feedback inhibition necessary to shut down activity and create an oscillation. Once the inhibition wears off it creates a window for the E cells to fire. The inhibitory time constants determine the spacing of these time-windows and thus the rhythmic frequency. Other relevant factors that determine the length of the oscillatory cycle is the input strength to the network, the pattern of connectivity between the E and I cells and the leak currents (Brunel and Wang, 2003). This mechanism has been termed “PING”, Pyramidal Interneuron Network Gamma, because it was originally conceived as a mechanism for gamma (Amit and Brunel, 1997; Brunel and Wang, 2003; Whittington et al., 2000). However, it can also generate beta-rhythmic ensembles (Lee et al., 2013; Lundqvist et al., 2011b). We should note that our hypotheses about the role of beta in WM do not depend on how beta is generated. We offer this as one possible mechanism; there are other models for beta generation (Sherman et al., 2016)

There could be two separate PING mechanisms in the superficial and in deep layers. Stronger gamma in superficial and stronger beta in deep layers resulting from different classes of I cells in superficial vs. deep layers with different time constants and/or the greater number of FS cells in superficial layers. The observed push-pull interaction between superficial gamma and deep beta could be generated by reciprocal inhibitory connections between the two PING networks (Lee et al., 2013).

The PING mechanism relies on strong excitation. For the gamma band, this drive is the sensory stimulation itself. In visual cortex, this generates strong, oscillatory gamma in response to sensory input which ceases when the stimulus is removed (Bastos et al., 2015a; Bosman et al., 2012; Brunet et al., 2013; Fries et al., 2008). In the PFC, gamma is more

bursty and variable, e.g., the center frequency varies and the bursts are sparse (Lundqvist et al., 2016). This may be because the PFC integrates inputs from many cortical and subcortical areas. Thus external sensory drive will have less of an impact on its overall excitation. Also, likely due to an enhanced number of excitatory connections on PFC cells (Elston, 2000), more lateral excitatory connectivity (Goldman-Rakic, 1996), longer intrinsic time constants (Murray et al., 2014) and synaptic mechanisms (Wang et al., 2006a), PFC networks are able to produce (bursty) gamma even in the absence of sensory stimuli.

The relationship between sensory input and rhythms is opposite for beta. Beta is more prominent in the absence of sensory drive (and, in somatomotor cortex, the absence of motor movement). Deep layer beta may be generated by a PING mechanism with excitatory drive provided via thalamocortical (Ketz et al., 2015) and/or basal ganglia (Chatham and Badre, 2015) loops that are self-sustaining in the absence of external inputs. Thus, beta is strong in the absence of sensory inputs, during planning, task set preparation, etc. Competition between the beta and gamma assemblies could control the “tuning” of the network to either internal (in beta) or external (in gamma) information (Brincat and Miller, 2016; Buschman and Miller, 2007).

Interplay between gamma and beta during WM control

To test whether this interplay between beta and gamma correlates with the control of WM, we used a sequence matching task (Lundqvist et al., 2018). Animals held sequences of two objects in WM and then had to judge whether a subsequent test sequence was a match. The advantage to this task is that it has multiple decision points. Animals have to determine if each object is a match both in identity and order. This affords more opportunity to examine WM control than a typically WM task which only involves remembering one stimulus and making one decision that co-occurs with a motor action.

This analysis revealed that shifts in the balance between beta and gamma/spiking did, in fact, correlate with WM control (Lundqvist et al., 2018). In anticipation of having to use a given object for the match decision (e.g., the first sample object for judging the first test object or the second for the second), there was reduced beta bursting along with an increase in gamma bursting and spiking information about the specific anticipated object. When an object held in WM was no longer needed, beta increased and gamma decreased together with spiking conveying information about that object. Further, deviations from “correct” beta/gamma dynamics predicted, not only a forthcoming error, but what kind of error the animal would make. For example, if the animal was going to mistakenly respond “match” to a non-matching sequence, the temporal dynamics of, and balance between, gamma and beta bursting looked like that on a match trial instead of non-match trial. We could also tell if the animals made the wrong decision to the first or second test object. In short, shifts in the balance between beta and gamma correlated with WM control processes; errors in the balance predicted upcoming behavioral errors.

Cortical Gradients

So far our discussion has emphasized the prefrontal cortex (PFC). Delay activity spiking, however, is a wide-spread cortical phenomenon (Fuster, 2015). But how widespread has recently generated vigorous debate, see (Christophel et al., 2017; Leavitt et al., 2017)). For example, Dotson and Gray recorded spiking activity from 42 cortical areas (Dotson et al., 2018). Delay activity was widespread but also showed gradients. In V1, delay activity was mostly decreased spiking in the delay relative to baseline, suggesting synaptic adaptation. This could also be a consequence of top-down signaling from higher cortical areas (van Kerkoerle et al., 2017). Other studies in early sensory areas also showed weak or non-existent delay activity spiking compared to higher-order cortical areas (Haller et al., 2018; Leavitt et al., 2017; Mendoza-Halliday et al., 2014). Interestingly, at the other extreme of cortical processing, in motor cortex, there is also little delay activity (Dotson et al., 2018; Haller et al., 2018). In between there is higher-order association cortex (including PFC, posterior parietal cortex, and temporal cortex) rich in delay activity (Dotson et al., 2018; Fuster, 1990; Haller et al., 2018; Leavitt et al., 2017; Sigala, 2009; Woloszyn and Sheinberg, 2009). These areas are highly interconnected (Markov et al., 2013b). They are also the cortical areas where top-down and bottom-up information reaches apex (Brincat et al., 2018; Siegel et al., 2015) and thus could support domain-general cognitive operations (Haller et al., 2018).

These higher-order delay activity-rich areas share several aspects of laminar circuitry. They have a balance in the soma size and cortico-cortical output connectivity between superficial vs. deep layers (Goulas et al., 2018). In contrast, motor output structures have a large laminar asymmetry in soma size (larger deep layer neurons) and predominant layer of cortical output (deep layers). Low-level sensory cortex features a highly differentiated and dense laminar circuit, emphasizing superficial layer soma size, and most corticocortical outputs originate from superficial layers.

We hypothesize that low-level sensory and motor cortex is not ideal for WM representation and control as a result of their local circuitry. Sensory areas have a relative emphasis on the superficial-layers (Zaldivar et al., 2018) where inputs can be richly encoded with gamma but lack the control element from deep layers. Motor areas emphasize deep layers (along with a predominance of beta), where outputs to motor structures can be gated but have a relatively poor superficial layer circuitry (Goulas et al., 2018). Association cortices lie in between. They have a relative balance between superficial and deep circuitry (Goulas et al., 2018) better suited for both representation and control of activity. In addition, there are other neuroanatomical gradients that also change from early to higher-order cortex, such as spine density and lateral connectivity (Elston, 2000; Goldman-Rakic, 1996). Both increase up the hierarchy making cells more intrinsically excitable and integrative (Murray et al., 2014; Wasmuht et al., 2018). The relative balance between specific inhibitory cell populations also changes (Kim et al., 2017), and could impact circuits for WM (Wang and Yang, 2018). It will be interesting to explore, in further work, which exact circuit elements enables higher-order cortex to sustain WM.

Putting it all together: A model for volitional control of working memory

Figure 3 illustrates our model. It shares many aspects with previously proposed circuits for visual sensory function (Bastos et al., 2015b; Mejias et al., 2016). Spikes encode and help maintain information in WM. Top-down information is associated with beta in deep (feedback) cortical layers (red wave). Bottom-up information is associated with gamma in superficial (feedforward) layers (blue waves). The central idea is that (top-down) deep-layer beta regulates the expression of (bottom-up) gamma in superficial layers thus gating the access of sensory information to WM and controlling its maintenance.

Alpha/beta and gamma oscillations can be below the threshold for spiking but they drive membrane potentials toward and sometimes over spike thresholds which is why there tends to be more spiking on the depolarizing phases of oscillations (Siegel et al., 2009)

Both superficial and deep layers of cortex are comprised of networks of deeply interconnected excitatory pyramidal (black) neurons and inhibitory (red) interneurons. Deep-layer beta is unidirectionally coupled to superficial layer beta. In turn, superficial-layer beta suppresses superficial-layer gamma oscillations. Note that the middle and deep layers of PFC are reciprocally connected with the mediodorsal nucleus of the thalamus, with layer 4 receiving thalamic input and layers 5 and 6 sending output to the thalamus (Giguere and Goldman-Rakic, 1988). WM delay interval spiking is prominent in MD thalamus (Watanabe and Funahashi, 2004). Beta-band coherence has been reported between PFC and thalamus during WM maintenance (Parnaudeau et al., 2013). Optogenetics suppression of MD thalamus suppresses cortical delay activity (Schmitt et al., 2017). Thus, the modulatory role of beta in the deep layers for WM control might be in part regulated by the thalamocortical loop.

To encode information in WM, deep-layer beta power and/or its coupling to superficial layer beta weakens. This disinhibits the recurrent excitation of layer 2/3 neurons (as indicated by the loop arrow) generating bursts of gamma. The gamma allows expression of spiking carrying bottom-up sensory inputs. The balance between beta and gamma can regulate the level of gamma bursting in the memory delay needed occasionally refresh the synaptic weight changes that help maintain the WMs. During WM read-out, beta is once again relaxed, allowing the increased gamma bursting and the ramp-up of spiking often seen near the end of memory delays (Hussar and Pasternak, 2010; Roesch and Olson, 2005). Increased spiking is needed so that WMs can acquire control of behavior. Keeping gamma bursting and spiking at a lower level earlier in the delay interval may prevent WMs from prematurely acquiring that control. To clear out WM, beta power/coupling increases. This suppresses gamma and the spiking that was maintaining the WM.

Summary and (many) open questions

Recent studies continue to indicate that memory delay spiking plays a critical mechanism for maintaining information in WM. But they also indicate that there is more going on than a simple persistence of spiking. Instead, there are brief bursts of spiking and associated gamma bursting that reflect activation and reactivation of the attractor states of the neural

ensembles for the WM memoranda. The spiking could cause temporary changes in synaptic weights that carry the WMs between spiking. This combination of spiking and synaptic weight changes solves many of the problems with persistent spiking. It is metabolically less expensive and makes the memories more robust to interference. It allows multiple items to be held in WM by “juggling” their activations in time. This new perspective is part of mounting evidence that the neural basis of cognition is not continuous (Buschman and Miller, 2010; Fiebelkorn et al., 2018; Helfrich et al., 2018; Landau and Fries, 2012; Schroeder et al., 2010; VanRullen, 2016). Sparse spiking also leaves room for rhythmic interplay between oscillations of different bands, gamma and alpha/beta. Beta is associated with top-down information and seems to have an inhibitory role. Increasing beta decreases gamma/spiking and vice-versa. Thus, the push-pull relationship (when beta is up, gamma is down and vice-versa) may be the infrastructure for top-down, executive control of WM storage. In short, beta can turn on and off the “faucet” of gamma-related WM reactivations.

Our discussion has been focused on WM and the higher cortical areas associated with WM. But we have noted that there is a similar laminar distribution of gamma and beta as well as a similar push-pull relationship between them all over sensory and motor cortex. Thus, rather than playing a role in WM only, this laminar interplay may be a cortical motif, a general mechanism by which the cortex, writ large, can control the inflow and processing of bottom-up sensory inputs via top-down knowledge (Bastos et al., 2012). This indeed fits with reports of widely distributed delay activity (Dotson et al., 2018; Fuster, 2015) and the idea that PFC could have more of a control function (Lara and Wallis, 2015; Miller and Cohen, 2001).

To be sure, this is just a beginning. Thus far, we have focused on the relatively simple processes of gating access to, and clearing out of, WM. But WM control is more than encoding and clearing information. It also involves manipulation. Information in WM can be transformed, reordered and sequenced, etc. This requires control at the level of individual ensembles, not just a general gating mechanism. Long-term memories can be loaded into WM; we do not yet know whether the same rhythmic interplay underlies this. What we are proposing is the infrastructure by which volition acts. There is also, of course, the big question of the genesis of volition itself. Our model, as with any other (including the classic model of WM storage) is just a starting point for more hypothesizing, further testing, and, if it is merited, updating.

Acknowledgements

This work was supported by NIMH R37MH087027, ONR MURI N00014-16-1-2832, The Picower Fellows Program, and the MIT Picower Institute Innovation Fund. We thank S.L. Brincat, M. Halassa, P. Herman, D. Pinotsis, K.B. Powell and A. Wutz for helpful discussions.

References

- Almeida R, Barbosa J, and Compte A (2015). Neural circuit basis of visuo-spatial working memory precision: a computational and behavioral study. *J Neurophysiol* 114, 1806–1818. [PubMed: 26180122]
- Amit DJ, and Brunel N (1997). Model of global spontaneous activity and local structured activity during delay periods in the cerebral cortex. *Cereb Cortex* 7, 237–252. [PubMed: 9143444]

- Antzoulatos EG, and Miller EK (2016). Synchronous beta rhythms of frontoparietal networks support only behaviorally relevant representations. *Elife* 5.
- Awh E, Barton B, and Vogel EK (2007). Visual Working Memory Represents a Fixed Number of Items Regardless of Complexity. *Visual Working Memory Represents a Fixed Number of Items Regardless of Complexity. Psychol Sci* 18, 622–628. [PubMed: 17614871]
- Bahramisharif A, Jensen O, Jacobs J, and Lisman J (2017). Serial representation of items during working memory maintenance at letter-selective cortical sites. *BioRxiv* 171660.
- Basar-Eroglu C, Brand A, Hildebrandt H, Karolina Kedzior K, Mathes B, and Schmiedt C (2007). Working memory related gamma oscillations in schizophrenia patients. *International Journal of Psychophysiology* 64, 39–45. [PubMed: 16962192]
- Bastos AM, Usrey WM, Adams RA, Mangun GR, Fries P, and Friston KJ (2012). Canonical microcircuits for predictive coding. *Neuron* 76, 695–711. [PubMed: 23177956]
- Bastos AM, Vezoli J, Bosman CA, Schoffelen J-M, Oostenveld R, Dowdall JR, De Weerd P, Kennedy H, and Fries P (2015a). Visual areas exert feedforward and feedback influences through distinct frequency channels. *Neuron* 85, 390–401. [PubMed: 25556836]
- Bastos AM, Litvak V, Moran R, Bosman CA, Fries P, and Friston KJ (2015b). A DCM study of spectral asymmetries in feedforward and feedback connections between visual areas V1 and V4 in the monkey. *Neuroimage* 108, 460–475. [PubMed: 25585017]
- Bastos AM, Loonis R, Kornblith S, Lundqvist M, and Miller EK (2018). Laminar recordings in frontal cortex suggest distinct layers for maintenance and control of working memory. *PNAS* 201710323.
- Bauer M, Oostenveld R, Peeters M, and Fries P (2006). Tactile Spatial Attention Enhances Gamma-Band Activity in Somatosensory Cortex and Reduces Low-Frequency Activity in Parieto-Occipital Areas. *J. Neurosci.* 26, 490–501. [PubMed: 16407546]
- Bollimunta A, Chen Y, Schroeder CE, and Ding M (2008). Neuronal mechanisms of cortical alpha oscillations in awake-behaving macaques. *The Journal of Neuroscience* 28, 9976. [PubMed: 18829955]
- Bonnefond M, and Jensen O (2012). Alpha oscillations serve to protect working memory maintenance against anticipated distracters. *Curr. Biol.* 22, 1969–1974. [PubMed: 23041197]
- Bosman CA, Schoffelen J-M, Brunet N, Oostenveld R, Bastos AM, Womelsdorf T, Rubehn B, Stieglitz T, De Weerd P, and Fries P (2012). Attentional Stimulus Selection through Selective Synchronization between Monkey Visual Areas. *Neuron* 75, 875–888. [PubMed: 22958827]
- Bressler SL, and Seth AK (2011). Wiener-Granger causality: a well established methodology. *Neuroimage* 58, 323–329. [PubMed: 20202481]
- Brincat SL, and Miller EK (2016). Prefrontal Cortex Networks Shift from External to Internal Modes during Learning. *J. Neurosci.* 36, 9739–9754. [PubMed: 27629722]
- Brincat SL, Siegel M, von Nicolai C, and Miller EK (2018). Gradual progression from sensory to task-related processing in cerebral cortex. *Proc. Natl. Acad. Sci. U.S.A*
- Brunel N, and Wang X-J (2003). What Determines the Frequency of Fast Network Oscillations With Irregular Neural Discharges? I. Synaptic Dynamics and Excitation-Inhibition Balance. *Journal of Neurophysiology* 90, 415–430. [PubMed: 12611969]
- Brunet N, Bosman CA, Roberts M, Oostenveld R, Womelsdorf T, De Weerd P, and Fries P (2013). Visual Cortical Gamma-Band Activity During Free Viewing of Natural Images. *Cereb. Cortex* 111, 3626–3631.
- Buffalo EA, Fries P, Landman R, Buschman TJ, and Desimone R (2011). Laminar differences in gamma and alpha coherence in the ventral stream. *Proceedings of the National Academy of Sciences* 108, 11262.
- Buschman TJ, and Miller EK (2007). Top-Down Versus Bottom-Up Control of Attention in the Prefrontal and Posterior Parietal Cortices. *Science* 315, 1860–1862. [PubMed: 17395832]
- Buschman TJ, and Miller EK (2010). Shifting the Spotlight of Attention: Evidence for Discrete Computations in Cognition. *Front Hum Neurosci* 4.
- Buschman TJ, Siegel M, Roy JE, and Miller EK (2011). Neural substrates of cognitive capacity limitations. *PNAS* 108, 11252–11255. [PubMed: 21690375]
- Buschman TJ, Denovellis EL, Diogo C, Bullock D, and Miller EK (2012). Synchronous oscillatory neural ensembles for rules in the prefrontal cortex. *Neuron* 76, 838–846. [PubMed: 23177967]

- Buzsáki G, and Wang X-J (2012). Mechanisms of gamma oscillations. *Annu. Rev. Neurosci* 35, 203–225. [PubMed: 22443509]
- Canolty RT, Edwards E, Dalal SS, Soltani M, Nagarajan SS, Kirsch HE, Berger MS, Barbaro NM, and Knight RT (2006). High gamma power is phase-locked to theta oscillations in human neocortex. *Science* 313, 1626–1628. [PubMed: 16973878]
- Chatham CH, and Badre D (2015). Multiple gates on working memory. *Current Opinion in Behavioral Sciences* 1, 23–31. [PubMed: 26719851]
- Christophel TB, Klink PC, Spitzer B, Roelfsema PR, and Haynes J-D (2017). The Distributed Nature of Working Memory. *Trends in Cognitive Sciences* 21, 111–124. [PubMed: 28063661]
- Colgin LL, Denninger T, Fyhn M, Hafting T, Bonnevie T, Jensen O, Moser M-B, and Moser EI (2009). Frequency of gamma oscillations routes flow of information in the hippocampus. *Nature* 462, 353–357. [PubMed: 19924214]
- Compte A, Brunel N, Goldman-Rakic PS, and Wang X-J (2000). Synaptic Mechanisms and Network Dynamics Underlying Spatial Working Memory in a Cortical Network Model. *Cereb Cortex* 10, 910–923. [PubMed: 10982751]
- Constantinidis C, Franowicz MN, and Goldman-Rakic PS (2001). Coding Specificity in Cortical Microcircuits: A Multiple-Electrode Analysis of Primate Prefrontal Cortex. *J. Neurosci.* 21, 3646–3655. [PubMed: 11331394]
- Cowan N (2010). The Magical Mystery Four The Magical Mystery Four: How Is Working Memory Capacity Limited, and Why? , How Is Working Memory Capacity Limited, and Why? *Curr Dir Psychol Sci* 19, 51–57. [PubMed: 20445769]
- Cromer JA, Roy JE, and Miller EK (2010). Representation of multiple, independent categories in the primate prefrontal cortex. *Neuron* 66, 796–807. [PubMed: 20547135]
- Dann B, Michaels JA, Schaffelhofer S, and Scherberger H (2016). Uniting functional network topology and oscillations in the fronto-parietal single unit network of behaving primates. *Elife* 5.
- David F, Courtiol E, Buonviso N, and Fourcaud-Trocmé N (2015). Competing Mechanisms of Gamma and Beta Oscillations in the Olfactory Bulb Based on Multimodal Inhibition of Mitral Cells Over a Respiratory Cycle. *ENeuro* 2.
- Desimone R, and Duncan J (1995). Neural mechanisms of selective visual attention. *Annu. Rev. Neurosci* 18, 193–222. [PubMed: 7605061]
- Donoghue JP, Sanes JN, Hatsopoulos NG, and Gaál G (1998). Neural discharge and local field potential oscillations in primate motor cortex during voluntary movements. *J. Neurophysiol* 79, 159–173. [PubMed: 9425187]
- Dotson NM, Hoffman SJ, Goodell B, and Gray CM (2018). Feature-Based Visual Short-Term Memory Is Widely Distributed and Hierarchically Organized. *Neuron*.
- Durstewitz D, Seamans JK, and Sejnowski TJ (2000). Neurocomputational models of working memory. *Nature Neuroscience* 3, 1184–1191. [PubMed: 11127836]
- Ede F van, Lange F de, Jensen O, and Maris E (2011). Orienting Attention to an Upcoming Tactile Event Involves a Spatially and Temporally Specific Modulation of Sensorimotor Alpha- and Beta-Band Oscillations. *J. Neurosci.* 31, 2016–2024. [PubMed: 21307240]
- van Ede F, Jensen O, and Maris E (2017). Supramodal Theta, Gamma, and Sustained Fields Predict Modality-specific Modulations of Alpha and Beta Oscillations during Visual and Tactile Working Memory. *J Cogn Neurosci* 29, 1455–1472. [PubMed: 28358658]
- Edin F, Klingberg T, Johansson P, McNab F, Tegnér J, and Compte A (2009). Mechanism for top-down control of working memory capacity. *PNAS* 106, 6802–6807. [PubMed: 19339493]
- Elston GN (2000). Pyramidal cells of the frontal lobe: all the more spinous to think with. *J. Neurosci.* 20, RC95. [PubMed: 10974092]
- Engle RW, Tuholski SW, Laughlin JE, and Conway AR (1999). Working memory, short-term memory, and general fluid intelligence: a latent-variable approach. *J Exp Psychol Gen* 128, 309–331. [PubMed: 10513398]
- Feingold J, Gibson DJ, DePasquale B, and Graybiel AM (2015). Bursts of beta oscillation differentiate postperformance activity in the striatum and motor cortex of monkeys performing movement tasks. *Proc. Natl. Acad. Sci. U.S.A* 112, 13687–13692. [PubMed: 26460033]

- Felleman DJ, and Van Essen DC (1991). Distributed hierarchical processing in the primate cerebral cortex. *Cerebral Cortex* 1, 1–47. [PubMed: 1822724]
- Fiebelkorn IC, Pinsk MA, and Kastner S (2018). A Dynamic Interplay within the Frontoparietal Network Underlies Rhythmic Spatial Attention. *Neuron* 99, 842–853.e8. [PubMed: 30138590]
- Fisch L, Privman E, Ramot M, Harel M, Nir Y, Kipervasser S, Andelman F, Neufeld MY, Kramer U, Fried I, et al. (2009). Neural “ignition”: enhanced activation linked to perceptual awareness in human ventral stream visual cortex. *Neuron* 64, 562–574. [PubMed: 19945397]
- Fontolan L, Morillon B, Liegeois-Chauvel C, and Giraud A-L (2014). The contribution of frequency-specific activity to hierarchical information processing in the human auditory cortex. *Nat Commun* 5, 4694. [PubMed: 25178489]
- Fries P (2015). Rhythms for Cognition: Communication through Coherence. *Neuron* 88, 220–235. [PubMed: 26447583]
- Fries P, Reynolds JH, Rorie AE, and Desimone R (2001). Modulation of oscillatory neuronal synchronization by selective visual attention. *Science* 291, 1560–1563. [PubMed: 11222864]
- Fries P, Womelsdorf T, Oostenveld R, and Desimone R (2008). The effects of visual stimulation and selective visual attention on rhythmic neuronal synchronization in macaque area V4. *J. Neurosci.* 28, 4823–4835. [PubMed: 18448659]
- Fuentemilla L, Penny WD, Cashdollar N, Bunzeck N, and Düzel E (2010). Theta-coupled periodic replay in working memory. *Curr. Biol* 20, 606–612. [PubMed: 20303266]
- Fujisawa S, Amarasingham A, Harrison MT, and Buzsáki G (2008). Behavior-dependent short-term assembly dynamics in the medial prefrontal cortex. *Nat. Neurosci* 11, 823–833. [PubMed: 18516033]
- Funahashi S, Bruce CJ, and Goldman-Rakic PS (1989). Mnemonic coding of visual space in the monkey’s dorsolateral prefrontal cortex. *Journal of Neurophysiology* 61, 331–349. [PubMed: 2918358]
- Fusi S, Miller EK, and Rigotti M (2016). Why neurons mix: high dimensionality for higher cognition. *Current Opinion in Neurobiology* 37, 66–74. [PubMed: 26851755]
- Fuster J (2015). *The Prefrontal Cortex* (Academic Press).
- Fuster JM (1990). Inferotemporal units in selective visual attention and short-term memory. *Journal of Neurophysiology* 64, 681–697. [PubMed: 2230917]
- Fuster JM (1999). *Memory in the Cerebral Cortex: An Empirical Approach to Neural Networks in the Human and Nonhuman Primate* (MIT Press).
- Fuster JM, and Alexander GE (1971a). Neuron activity related to short-term memory. *Science* 173, 652–654. [PubMed: 4998337]
- Fuster JM, and Alexander GE (1971b). Neuron activity related to short-term memory. *Science* 173, 652–654. [PubMed: 4998337]
- Giguere M, and Goldman-Rakic PS (1988). Mediodorsal nucleus: areal, laminar, and tangential distribution of afferents and efferents in the frontal lobe of rhesus monkeys. *J. Comp. Neurol* 277, 195–213. [PubMed: 2466057]
- Gilbert CD, and Wiesel TN (1983). Functional organization of the visual cortex. *Prog. Brain Res* 58, 209–218. [PubMed: 6138809]
- Goldman-Rakic P. (1995). Cellular basis of working memory. *Neuron* 14, 477–485. [PubMed: 7695894]
- Goldman-Rakic PS (1996). Regional and cellular fractionation of working memory. *Proc. Natl. Acad. Sci. U.S.A* 93, 13473–13480. [PubMed: 8942959]
- Goulas A, Zilles K, and Hilgetag CC (2018). Cortical Gradients and Laminar Projections in Mammals. *Trends Neurosci.*
- Gregoriou GG, Gotts SJ, Zhou H, and Desimone R (2009). High-Frequency, Long-Range Coupling Between Prefrontal and Visual Cortex During Attention. *Science* 324, 1207–1210. [PubMed: 19478185]
- Haegens S, Nacher V, Luna R, Romo R, and Jensen O (2011). α -Oscillations in the monkey sensorimotor network influence discrimination performance by rhythmical inhibition of neuronal spiking. *Proc. Natl. Acad. Sci. U.S.A* 108, 19377–19382. [PubMed: 22084106]

- Haller M, Case J, Crone NE, Chang EF, King-Stephens D, Laxer KD, Weber PB, Parvizi J, Knight RT, and Shestyuk AY (2018). Persistent neuronal activity in human prefrontal cortex links perception and action. *Nat Hum Behav* 2, 80–91. [PubMed: 29963646]
- Helfrich RF, Fiebelkorn IC, Szczepanski SM, Lin JJ, Parvizi J, Knight RT, and Kastner S (2018). Neural Mechanisms of Sustained Attention Are Rhythmic. *Neuron* 99, 854–865.e5. [PubMed: 30138591]
- Herman PA, Lundqvist M, and Lansner A (2013). Nested theta to gamma oscillations and precise spatiotemporal firing during memory retrieval in a simulated attractor network. *Brain Res.* 1536, 68–87. [PubMed: 23939226]
- Honkanen R, Rouhinen S, Wang SH, Palva JM, and Palva S (2015). Gamma Oscillations Underlie the Maintenance of Feature-Specific Information and the Contents of Visual Working Memory. *Cereb. Cortex* 25, 3788–3801. [PubMed: 25405942]
- Howard MW, Rizzuto DS, Caplan JB, Madsen JR, Lisman J, Aschenbrenner-Scheibe R, Schulze-Bonhage A, and Kahana MJ (2003). Gamma oscillations correlate with working memory load in humans. *Cereb. Cortex* 13, 1369–1374. [PubMed: 14615302]
- Hussar C, and Pasternak T (2010). Trial-to-trial variability of the prefrontal neurons reveals the nature of their engagement in a motion discrimination task. *PNAS* 201009956.
- Hussar CR, and Pasternak T (2012). Memory-Guided Sensory Comparisons in the Prefrontal Cortex: Contribution of Putative Pyramidal Cells and Interneurons. *J. Neurosci* 32, 2747–2761. [PubMed: 22357858]
- Jensen O, and Mazaheri A (2010). Shaping Functional Architecture by Oscillatory Alpha Activity: Gating by Inhibition. *Frontiers in Human Neuroscience* 4.
- Jha A, Nachev P, Barnes G, Husain M, Brown P, and Litvak V (2015). The Frontal Control of Stopping. *Cereb Cortex* 25, 4392–4406. [PubMed: 25754518]
- Just MA, and Carpenter PA (1992). A capacity theory of comprehension: individual differences in working memory. *Psychol Rev* 99, 122–149. [PubMed: 1546114]
- van Kerkoerle T, Self MW, Dagnino B, Gariel-Mathis M-A, Poort J, van der Togt C, and Roelfsema PR (2014). Alpha and gamma oscillations characterize feedback and feedforward processing in monkey visual cortex. *Proc. Natl. Acad. Sci. U.S.A* 111, 14332–14341. [PubMed: 25205811]
- van Kerkoerle T, Self MW, and Roelfsema PR (2017). Layer-specificity in the effects of attention and working memory on activity in primary visual cortex. *Nat Commun* 8, 13804. [PubMed: 28054544]
- Ketz NA, Jensen O, and O'Reilly RC (2015). Thalamic pathways underlying prefrontal cortex-medial temporal lobe oscillatory interactions. *Trends Neurosci.* 38, 3–12. [PubMed: 25455705]
- Kim Y, Yang GR, Pradhan K, Venkataraju KU, Bota M, García Del Molino LC, Fitzgerald G, Ram K, He M, Levine JM, et al. (2017). Brain-wide Maps Reveal Stereotyped Cell-Type-Based Cortical Architecture and Subcortical Sexual Dimorphism. *Cell* 171, 456–469.e22. [PubMed: 28985566]
- Kopell N, Whittington MA, and Kramer MA (2011). Neuronal assembly dynamics in the beta1 frequency range permits short-term memory. *Proc. Natl. Acad. Sci. U.S.A* 108, 3779–3784. [PubMed: 21321198]
- Kornblith S, Buschman TJ, and Miller EK (2016). Stimulus Load and Oscillatory Activity in Higher Cortex. *Cereb Cortex* 26, 3772–3784. [PubMed: 26286916]
- Kritzer MF, and Goldman-Rakic PS (1995). Intrinsic circuit organization of the major layers and sublayers of the dorsolateral prefrontal cortex in the rhesus monkey. *Journal of Comparative Neurology* 359, 131–143.
- Kucewicz MT, Berry BM, Kremen V, Brinkmann BH, Sperling MR, Jobst BC, Gross RE, Lega B, Sheth SA, Stein JM, et al. (2017). Dissecting gamma frequency activity during human memory processing. *Brain* 140, 1337–1350. [PubMed: 28335018]
- Lakatos P, Shah AS, Knuth KH, Ulbert I, Karmos G, and Schroeder CE (2005). An oscillatory hierarchy controlling neuronal excitability and stimulus processing in the auditory cortex. *J. Neurophysiol* 94, 1904–1911. [PubMed: 15901760]
- Landau AN, and Fries P (2012). Attention samples stimuli rhythmically. *Curr. Biol* 22, 1000–1004. [PubMed: 22633805]

- Lara AH, and Wallis JD (2015). The Role of Prefrontal Cortex in Working Memory: A Mini Review. *Front. Syst. Neurosci* 9.
- Leavitt ML, Mendoza-Halliday D, and Martinez-Trujillo JC (2017). Sustained Activity Encoding Working Memories: Not Fully Distributed. *Trends Neurosci.* 40, 328–346. [PubMed: 28515011]
- Lee JH, Whittington MA, and Kopell NJ (2013). Top-Down Beta Rhythms Support Selective Attention via Interlaminar Interaction: A Model. *PLoS Comput. Biol* 9, e1003164. [PubMed: 23950699]
- Leszczynski M, Fell J, Jensen O, and Axmacher N (2017). Alpha activity in the ventral and dorsal visual stream controls information flow during working memory. *BioRxiv* 180166.
- Lisman J (1997). Bursts as a unit of neural information: making unreliable synapses reliable. *Trends in Neurosciences* 20, 38–43. [PubMed: 9004418]
- Lisman JE, and Idiart MA (1995). Storage of 7 ± 2 short-term memories in oscillatory subcycles. *Science* 267, 1512–1515. [PubMed: 7878473]
- Luck SJ, and Vogel EK (1997). The capacity of visual working memory for features and conjunctions. *Nature* 390, 279–281. [PubMed: 9384378]
- Lundqvist M, Herman P, and Lansner A (2011a). Theta and gamma power increases and alpha/beta power decreases with memory load in an attractor network model. *J Cogn Neurosci* 23, 3008–3020. [PubMed: 21452933]
- Lundqvist M, Herman P, and Lansner A (2011b). Theta and gamma power increases and alpha/beta power decreases with memory load in an attractor network model. *J Cogn Neurosci* 23, 3008–3020. [PubMed: 21452933]
- Lundqvist M, Herman P, and Lansner A (2012). Variability of spike firing during θ -coupled replay of memories in a simulated attractor network. *Brain Res* 1434, 152–161. [PubMed: 21907326]
- Lundqvist M, Rose J, Herman P, Brincat SL, Buschman TJ, and Miller EK (2016). Gamma and Beta Bursts Underlie Working Memory. *Neuron* 90, 152–164. [PubMed: 26996084]
- Lundqvist M, Herman P, Warden MR, Brincat SL, and Miller EK (2018). Gamma and beta bursts during working memory readout suggest roles in its volitional control. *Nat Commun* 9, 394. [PubMed: 29374153]
- Lundqvist M, Herman P, Miller EK, M. (in press). Working Memory: Delay Activity, Yes! Persistent Activity? Maybe Not. *Journal of Neuroscience*.
- Maier A, Adams GK, Aura C, and Leopold DA (2010). Distinct superficial and deep laminar domains of activity in the visual cortex during rest and stimulation. *Frontiers in Systems Neuroscience* 4.
- Markov NT, Vezoli J, Chameau P, Falchier A, Quilodran R, Huissoud C, Lamy C, Misery P, Giroud P, Ullman S, et al. (2013a). The anatomy of hierarchy: Feedforward and feedback pathways in macaque visual cortex. *J. Comp. Neurol.*
- Markov NT, Ercsey-Ravasz M, Van Essen DC, Knoblauch K, Toroczkai Z, and Kennedy H (2013b). Cortical high-density counterstream architectures. *Science* 342, 1238406. [PubMed: 24179228]
- Mejias JF, Murray JD, Kennedy H, and Wang X-J (2016). Feedforward and feedback frequency-dependent interactions in a large-scale laminar network of the primate cortex. *Science Advances* 2, e1601335–e1601335. [PubMed: 28138530]
- Mendoza-Halliday D, Torres S, and Martinez-Trujillo JC (2014). Sharp emergence of feature-selective sustained activity along the dorsal visual pathway. *Nat. Neurosci* 17, 1255–1262. [PubMed: 25108910]
- Meyers EM, Freedman DJ, Kreiman G, Miller EK, and Poggio T (2008). Dynamic Population Coding of Category Information in Inferior Temporal and Prefrontal Cortex. *Journal of Neurophysiology* 100, 1407–1419. [PubMed: 18562555]
- Mi Y, Katkov M, and Tsodyks M (2017). Synaptic Correlates of Working Memory Capacity. *Neuron* 93, 323–330. [PubMed: 28041884]
- Michalareas G, Vezoli J, van Pelt S, Schoffelen J-M, Kennedy H, and Fries P (2016). Alpha-Beta and Gamma Rhythms Subserve Feedback and Feedforward Influences among Human Visual Cortical Areas. *Neuron* 89, 384–397. [PubMed: 26777277]
- Miller EK, and Cohen JD (2001). An Integrative Theory of Prefrontal Cortex Function. *Annual Review of Neuroscience* 24, 167–202.

- Miller EK, and Wilson MA (2008). All my circuits: using multiple electrodes to understand functioning neural networks. *Neuron* 60, 483–488. [PubMed: 18995823]
- Miller EK, Erickson CA, and Desimone R (1996). Neural mechanisms of visual working memory in prefrontal cortex of the macaque. *Journal of Neuroscience* 16, 5154–5167. [PubMed: 8756444]
- Mongillo G, Barak O, and Tsodyks M (2008). Synaptic Theory of Working Memory. *Science* 319, 1543–1546. [PubMed: 18339943]
- Murray JD, Bernacchia A, Freedman DJ, Romo R, Wallis JD, Cai X, Padoa-Schioppa C, Pasternak T, Seo H, Lee D, et al. (2014). A hierarchy of intrinsic timescales across primate cortex. *Nat. Neurosci* 17, 1661–1663. [PubMed: 25383900]
- Murray JD, Bernacchia A, Roy NA, Constantinidis C, Romo R, and Wang X-J (2017). Stable population coding for working memory coexists with heterogeneous neural dynamics in prefrontal cortex. *PNAS* 114, 394–399. [PubMed: 28028221]
- Naud R, and Sprekeler H (2018). Sparse bursts optimize information transmission in a multiplexed neural code. *PNAS* 115, E6329–E6338. [PubMed: 29934400]
- Parnaudeau S, O’Neill P-K, Bolkan SS, Ward RD, Abbas AI, Roth BL, Balsam PD, Gordon JA, and Kellendonk C (2013). Inhibition of Mediodorsal Thalamus Disrupts Thalamofrontal Connectivity and Cognition. *Neuron* 77, 1151–1162. [PubMed: 23522049]
- Parthasarathy A, Herikstad R, Bong JH, Medina FS, Libedinsky C, and Yen S-C (2017). Mixed selectivity morphs population codes in prefrontal cortex. *Nature Neuroscience* 1.
- Pasternak T, and Greenlee MW (2005). Working memory in primate sensory systems. *Nature Reviews Neuroscience* 6, 97–107. [PubMed: 15654324]
- Popov T, Kastner S, and Jensen O (2017). FEF-Controlled Alpha Delay Activity Precedes Stimulus-Induced Gamma-Band Activity in Visual Cortex. *J. Neurosci* 37, 4117–4127. [PubMed: 28314817]
- Richter CG, Thompson WH, Bosman CA, and Fries P (2017). Top-Down Beta Enhances Bottom-Up Gamma. *J. Neurosci.* 37, 6698–6711. [PubMed: 28592697]
- Richter CG, Coppola R, and Bressler SL (2018). Top-down beta oscillatory signaling conveys behavioral context in early visual cortex. *Sci Rep* 8, 6991. [PubMed: 29725028]
- Rigotti M, Barak O, Warden MR, Wang X-J, Daw ND, Miller EK, and Fusi S (2013). The importance of mixed selectivity in complex cognitive tasks. *Nature* 497, 585–590. [PubMed: 23685452]
- Rockland KS, and Pandya DN (1979). Laminar origins and terminations of cortical connections of the occipital lobe in the rhesus monkey. *Brain Res.* 179, 3–20. [PubMed: 116716]
- Roesch MR, and Olson CR (2005). Neuronal Activity Dependent on Anticipated and Elapsed Delay in Macaque Prefrontal Cortex, Frontal and Supplementary Eye Fields, and Premotor Cortex. *Journal of Neurophysiology* 94, 1469–1497. [PubMed: 15817652]
- Romo R, Brody CD, Hernández A, and Lemus L (1999). Neuronal correlates of parametric working memory in the prefrontal cortex. *Nature* 399, 470–473. [PubMed: 10365959]
- Rose NS, LaRocque JJ, Riggall AC, Gosseries O, Starrett MJ, Meyering EE, and Postle BR (2016). Reactivation of latent working memories with transcranial magnetic stimulation. *Science* 354, 1136–1139. [PubMed: 27934762]
- Roux F, Wibral M, Mohr HM, Singer W, and Uhlhaas PJ (2012). Gamma-band activity in human prefrontal cortex codes for the number of relevant items maintained in working memory. *J. Neurosci.* 32, 12411–12420. [PubMed: 22956832]
- Salazar RF, Dotson NM, Bressler SL, and Gray CM (2012). Content-Specific Fronto-Parietal Synchronization During Visual Working Memory. *Science* 1224000.
- Sandberg A, Tegnér J, and Lansner A (2003). A working memory model based on fast Hebbian learning. *Network* 14, 789–802. [PubMed: 14653503]
- Schmitt LI, Wimmer RD, Nakajima M, Happ M, Mofakham S, and Halassa MM (2017). Thalamic amplification of cortical connectivity sustains attentional control. *Nature* 545, 219–223. [PubMed: 28467827]
- Schroeder CE, Wilson DA, Radman T, Scharfman H, and Lakatos P (2010). Dynamics of Active Sensing and perceptual selection. *Current Opinion in Neurobiology* 20, 172–176. [PubMed: 20307966]

- 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. [PubMed: 17418956]
- Sherman MA, Lee S, Law R, Haegens S, Thorn CA, Hämäläinen MS, Moore CI, and Jones SR (2016). Neural mechanisms of transient neocortical beta rhythms: Converging evidence from humans, computational modeling, monkeys, and mice. *Proc. Natl. Acad. Sci. U.S.A* 113, E4885–4894. [PubMed: 27469163]
- Siegel M, Warden MR, and Miller EK (2009). Phase-dependent neuronal coding of objects in short-term memory. *PNAS* 106, 21341–21346. [PubMed: 19926847]
- Siegel M, Buschman TJ, and Miller EK (2015). Cortical information flow during flexible sensorimotor decisions. *Science* 348, 1352–1355. [PubMed: 26089513]
- Sigala N (2009). Visual Working Memory and Delay Activity in Highly Selective Neurons in the Inferior Temporal Cortex. *Front Syst Neurosci* 3.
- Smith MA, Jia X, Zandvakili A, and Kohn A (2013). Laminar dependence of neuronal correlations in visual cortex. *J. Neurophysiol.* 109, 940–947. [PubMed: 23197461]
- Spaak E, Bonnefond M, Maier A, Leopold DA, and Jensen O (2012). Layer-specific entrainment of γ -band neural activity by the α rhythm in monkey visual cortex. *Curr. Biol.* 22, 2313–2318. [PubMed: 23159599]
- Spaak E, Watanabe K, Funahashi S, and Stokes MG (2017a). Stable and Dynamic Coding for Working Memory in Primate Prefrontal Cortex. *J. Neurosci.* 37, 6503–6516. [PubMed: 28559375]
- Spaak E, Watanabe K, Funahashi S, and Stokes MG (2017b). Stable and dynamic coding for working memory in primate prefrontal cortex. *J. Neurosci.* 3364–16. [PubMed: 28258168]
- Spitzer B, and Haegens S (2017). Beyond the Status Quo: A Role for Beta Oscillations in Endogenous Content (Re-) Activation. *ENeuro* ENEURO.0170-17.2017.
- Sprague TC, Ester EF, and Serences JT (2016). Restoring Latent Visual Working Memory Representations in Human Cortex. *Neuron* 91, 694–707. [PubMed: 27497224]
- Stanley DA, Roy JE, Aoi MC, Kopell NJ, and Miller EK (2018). Low-Beta Oscillations Turn Up the Gain During Category Judgments. *Cereb. Cortex* 28, 116–130. [PubMed: 29253255]
- Stokes MG (2015). ‘Activity-silent’ working memory in prefrontal cortex: a dynamic coding framework. *Trends in Cognitive Sciences* 19, 394–405. [PubMed: 26051384]
- Stokes M, and Spaak E (2016). The Importance of Single-Trial Analyses in Cognitive Neuroscience. *Trends in Cognitive Sciences* 20, 483–486. [PubMed: 27237797]
- Stokes MG, Kusunoki M, Sigala N, Nili H, Gaffan D, and Duncan J (2013). Dynamic Coding for Cognitive Control in Prefrontal Cortex. *Neuron* 78, 364–375. [PubMed: 23562541]
- Uhlhaas PJ, and Singer W (2010). Abnormal neural oscillations and synchrony in schizophrenia. *Nature Reviews Neuroscience* 11, 100. [PubMed: 20087360]
- VanRullen R (2016). Perceptual Cycles. *Trends in Cognitive Sciences* 20, 723–735. [PubMed: 27567317]
- Vogel EK, and Machizawa MG (2004). Neural activity predicts individual differences in visual working memory capacity. *Nature* 428, 748–751. [PubMed: 15085132]
- Voytek B, Kayser AS, Badre D, Fegen D, Chang EF, Crone NE, Parvizi J, Knight RT, and D’Esposito M (2015). Oscillatory dynamics coordinating human frontal networks in support of goal maintenance. *Nat. Neurosci.* 18, 1318–1324. [PubMed: 26214371]
- Wang X-J (2010). Neurophysiological and computational principles of cortical rhythms in cognition. *Physiol. Rev* 90, 1195–1268. [PubMed: 20664082]
- Wang X-J, and Yang GR (2018). A disinhibitory circuit motif and flexible information routing in the brain. *Current Opinion in Neurobiology* 49, 75–83. [PubMed: 29414069]
- Wang Y, Markram H, Goodman PH, Berger TK, Ma J, and Goldman-Rakic PS (2006). Heterogeneity in the pyramidal network of the medial prefrontal cortex. *Nat. Neurosci.* 9, 534–542. [PubMed: 16547512]
- Warden MR, and Miller EK (2010). Task-dependent changes in short-term memory in the prefrontal cortex. *J. Neurosci.* 30, 15801–15810. [PubMed: 21106819]

- Wasmuht DF, Spaak E, Buschman TJ, Miller EK, and Stokes MG (2018). Intrinsic neuronal dynamics predict distinct functional roles during working memory. *Nature Communications* 9, 3499.
- Watanabe K, and Funahashi S (2014). Neural mechanisms of dual-task interference and cognitive capacity limitation in the prefrontal cortex. *Nature Neuroscience* 17, 601–611. [PubMed: 24584049]
- Watanabe Y, and Funahashi S (2004). Neuronal activity throughout the primate mediodorsal nucleus of the thalamus during oculomotor delayed-responses. I. Cue-, delay-, and response-period activity. *J. Neurophysiol.* 92, 1738–1755. [PubMed: 15140911]
- Watrous AJ, Fell J, Ekstrom AD, and Axmacher N (2015). More than spikes: common oscillatory mechanisms for content specific neural representations during perception and memory. *Current Opinion in Neurobiology* 31, 33–39. [PubMed: 25129044]
- Whittington MA, Traub RD, Kopell N, Ermentrout B, and Buhl EH (2000). Inhibition-based rhythms: experimental and mathematical observations on network dynamics. *Int J Psychophysiol* 38, 315–336. [PubMed: 11102670]
- Wolff MJ, Jochim J, Akyürek EG, and Stokes MG (2017). Dynamic hidden states underlying working-memory-guided behavior. *Nat. Neurosci.* 20, 864–871. [PubMed: 28414333]
- Woloszyn L, and Sheinberg DL (2009). Neural Dynamics in Inferior Temporal Cortex during a Visual Working Memory Task. *J. Neurosci.* 29, 5494–5507. [PubMed: 19403817]
- Wutz A, Loonis R, Roy JE, Donoghue JA, and Miller EK (2018). Different Levels of Category Abstraction by Different Dynamics in Different Prefrontal Areas. *Neuron* 97, 716–726.e8. [PubMed: 29395915]
- Xing D, Yeh C-I, Burns S, and Shapley RM (2012). Laminar analysis of visually evoked activity in the primary visual cortex. *Proc. Natl. Acad. Sci. U.S.A* 109, 13871–13876. [PubMed: 22872866]
- Zaldivar D, Goense J, Lowe SC, Logothetis NK, and Panzeri S (2018). Dopamine Is Signaled by Midfrequency Oscillations and Boosts Output Layers Visual Information in Visual Cortex. *Current Biology* 28, 224–235.e5. [PubMed: 29307559]
- Zhang Y, Chen Y, Bressler SL, and Ding M (2008). Response preparation and inhibition: The role of the cortical sensorimotor beta rhythm. *Neuroscience* 156, 238–246. [PubMed: 18674598]

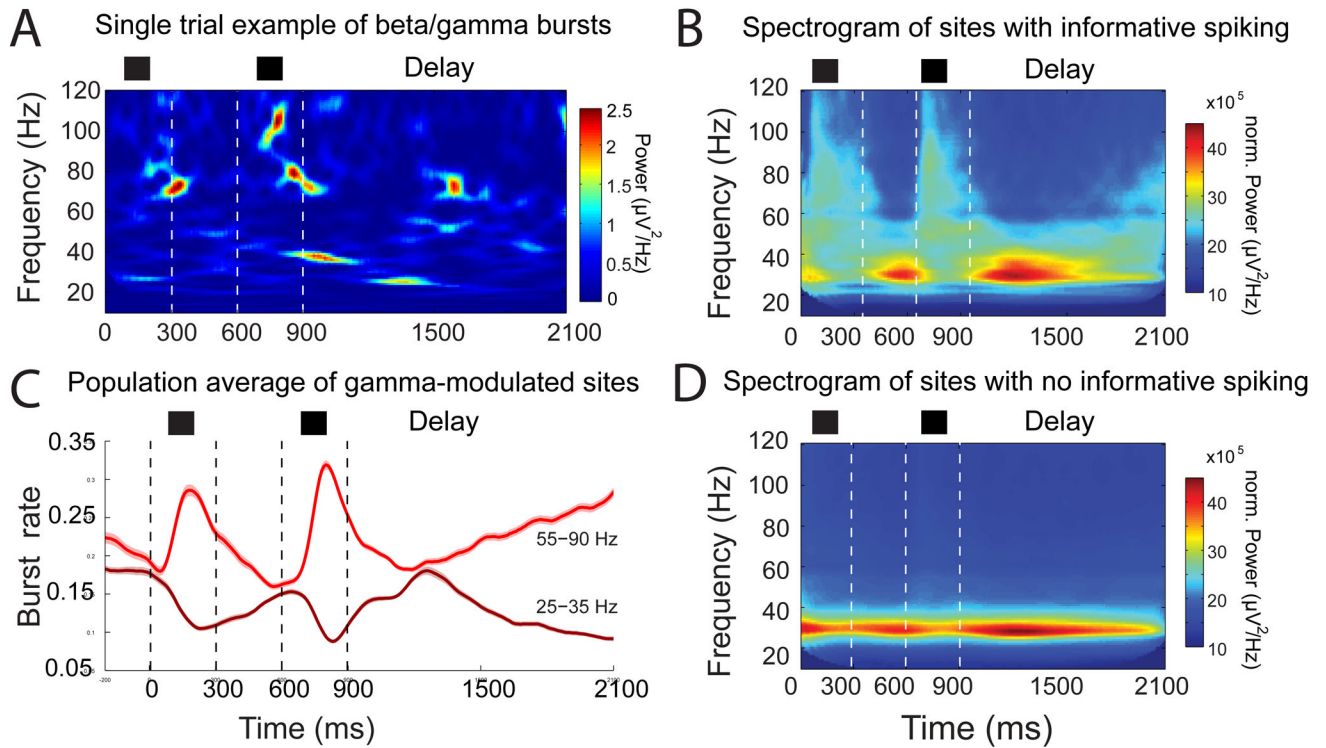


Figure 1 – Gamma and beta bursts underlie working memory.

A) A single trial example of LFP power in time and frequency. Two stimuli were presented (S1 and S2) and later tested following a delay. Narrow bursts of power in the beta and gamma bands are evident both during cue processing and delay. LFP data from sites that contained spikes which carried information about the presented cue (B) vs. those that did not (D) are shown. Only sites containing informative spiking (D, population average) showed modulation of beta and gamma. This effect remained after controlling for differences in spike rate between informative (B) vs. non-informative (C) sites. C) On gamma-modulated sites, the beta and gamma burst rates are mirror images of each other. Gamma bursting increases during stimulus presentation and towards the end of the delay, and beta does the opposite. D) On sites without informative spiking, only beta is task modulated and less so. Modified from Lundqvist et al (2016)

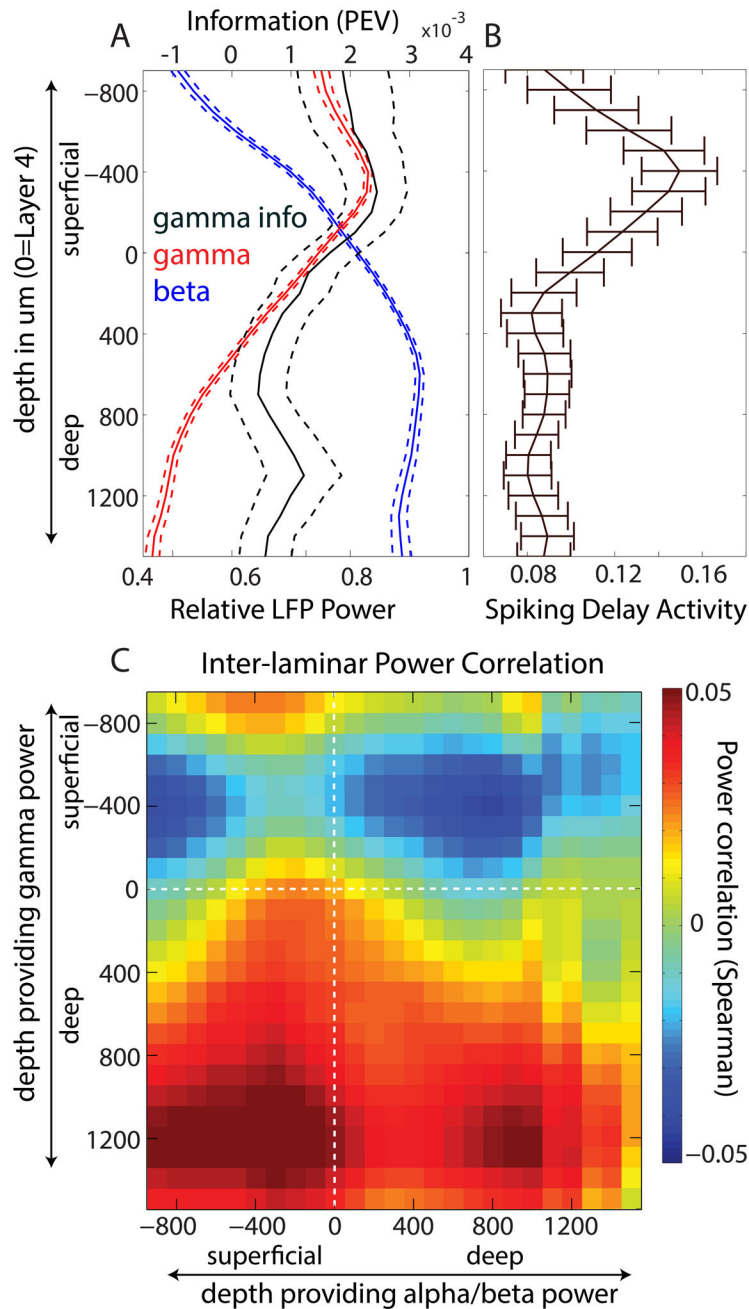


Figure 2 –. Laminar organization of gamma/beta rhythms and delay activity.

A) Gamma and alpha/beta power are segregated into distinct layers. Gamma power peaks 400 um above layer 4 whereas alpha/beta power peaks at 600 um below layer 4. Gamma bursts in superficial, but not deep layers, carry significant information about the cued item during the WM delay period (quantified by the Percent Explained Variance, or PEV, statistic). Beta bursts do not carry significant information during the delay (not shown). Dotted lines are ± 1 SEM across sessions (N=60). B) Spiking activity, quantified by multiunit change from baseline (arbitrary units) during the delay period is strongest in superficial layers. The pattern of laminar pattern of delay activity correlates strongly with

gamma, and is strongly anti-correlated with alpha/beta. C) Correlation map between gamma and beta power across layers. Deep layer beta power is anti-correlated with superficial layer gamma power during the WM delay. From Bastos et al (2018).

Author Manuscript

Author Manuscript

Author Manuscript

Author Manuscript

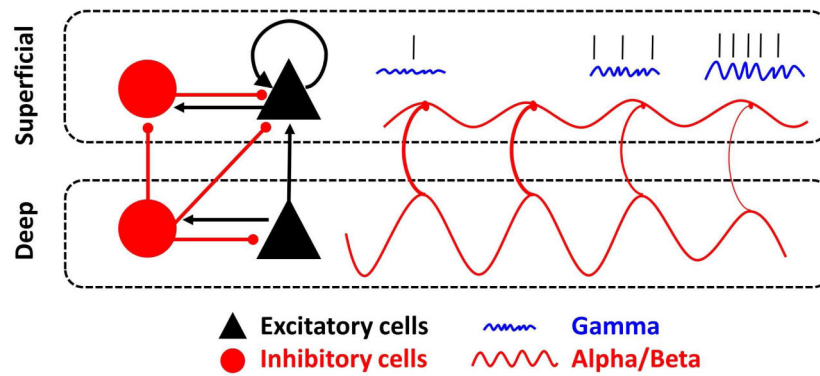


Figure 3 –. A model of WM.

Denoted by two rectangular, dashed boxes, two cortical compartments, superficial and deep, are made up of densely interconnected excitatory pyramidal (black) and inhibitory (red) interneurons. Inhibitory connections are line segments with a red, rounded end, and excitatory connections are line segments with a black, arrow end. Two separate PING networks in superficial vs. deep layers are responsible for generating gamma in superficial layers and beta in deep layers (sustained by connections to thalamus and basal ganglia, not shown). The looping arrow returning on itself in the superficial layers represents the recurrent connectivity found within layer 3 pyramidal cell networks in prefrontal cortex. The sinusoidal red-line in deep layers reflects beta oscillations and their driving influence on superficial beta oscillations. Beta oscillations are phase-amplitude coupled with gamma oscillations (blue squiggly lines), and these gamma oscillations organize delay-period spiking representing WM content (straight black marks). Spiking activity inside gamma bursts is more informative than outside. Over time, moving from left to right in the figure, the deep beta reduces in power and releases inhibition onto the superficial layers. This results in enhanced superficial gamma and spiking, i.e., enhanced maintenance of WM, as is seen when transitioning between baseline to WM task performance. The reversed process (enhancement of deep layer beta, enhanced suppression of superficial layer gamma/spiking) which would “clear out” the contents of WM, as seen at the end of the trial, or when WM contents are no longer needed. From Bastos et al (2018)