

Published in final edited form as:

Trends Cogn Sci. 2017 June 01; 21(6): 449–461. doi:10.1016/j.tics.2017.03.010.

Prioritizing Information during Working Memory: Beyond Sustained Internal Attention

Nicholas E. Myers, Mark G. Stokes, Anna C. Nobre

Oxford Centre for Human Brain Activity and Department of Experimental Psychology, University of Oxford

Abstract

Working memory (WM) capacity limits give attention the important mandate of gating in only relevant information. It is increasingly evident that attention is equally crucial for prioritizing representations within WM as the importance of individual items changes. Retrospective prioritization has been proposed to result from a focus of internal attention highlighting one of several representations. We suggest an updated model, in which prioritization acts in multiple steps: first orienting towards and selecting a memory, and then reconfiguring its representational state in the service of upcoming task demands. Reconfiguration sets up an optimized perception-action mapping, obviating the need for sustained attention. This view is consistent with recent literature, makes testable predictions, and links WM with task-switching and action preparation.

Keywords

Working Memory; Attention; Focus of Attention; Retrocue; Task Set; Latent Information Storage

The Changing Concept of Priority in Working Memory

The subject of this review is the neural basis and behavioral consequence of prioritizing information maintained in visual short-term, ‘working’ memory (WM). By working memory in this context we refer to the ability to store and manipulate recently acquired information for a period of seconds, independently of continuous sensory stimulation, to guide behavior over the short-term¹. This ability is central to intelligent behavior [1,2], and therefore touches on nearly all domains of cognitive neuroscience (such as fluid intelligence, perceptual decision-making, or model-based learning, see e.g., [3–5]). The severe limits on how much can be encoded in working memory – conceived as a small number of quantized representations [6–8] or as a limited pool of mnemonic resources [9] – hamper our ability to act optimally when there is too much information to be considered at once. As a consequence of this bottleneck, attention is of central importance to working memory [10–15]: Those who cannot select the most important information and keep out irrelevant distraction unnecessarily clutter their working memory store [16,17].

correspondence to: kia.nobre@ohba.ox.ac.uk or nicholas.myers@ohba.ox.ac.uk.

¹This definition is independent of the psychological quality of maintaining the memoranda in awareness (in mind).

Early studies exploring the role of attention in WM manipulated selective encoding (i.e., prioritizing a subset of items during encoding). Later, studies revealed that focusing on the relevant pieces of information even after they have already been encoded also improves memory [18–20]. Such retrospective cueing cannot influence basic sensory processing of the memory items, or their encoding, but rather operates at a pure mnemonic level, prioritizing the contents maintained in WM.

Neurodevelopmental [e.g., 21] and psychiatric disorders [e.g., 22], as well as healthy ageing [23,24], severely affect working memory capacity, making it imperative that we better understand how prioritization within WM can help us make the most of a preciously limited cognitive resource. This review focuses on new empirical and theoretical advances that shed light on the role of prioritization in working memory, and how this may relate to task preparation. In synthesizing this literature, we suggest that both attentional selection and task preparation play critical roles in prioritizing information in working memory to guide optimal performance.

We begin by drawing parallels with the better-understood mechanisms of selective attention for perception. We then build on this model with the aim of explaining more fully how prioritization may operate in WM, and within internal information stores more generally. We propose that, in addition to any benefits brought about by attentional selection of individual items, behavioral benefits also arise in large part because of preparation of the right behavioral policy (for instance, by setting up appropriate contingencies between upcoming stimuli and behavior). Our proposal can account for a number of otherwise odd findings in the behavioral literature. Moreover, it may help pin down the dual roles of selection and preparation in prioritizing information in mind. Furthermore, our model makes predictions about the possible neural basis of the architecture of WM.

Attention in Perception and Working Memory

WM is famously burdened with severe capacity limits. As in many other domains of cognition that contain a bottleneck [25,26], the preferential selection of pertinent information seems crucial if we are to make the best use of our limited resource. In the domain of perception, the term ‘selective attention’ is invoked to describe such preferential biases towards behaviorally relevant stimuli. In extending this literature, attention has been shown to be influential for selecting information for encoding into WM [27–29], and for preventing distracting information from gaining access to it [30]. The benefits of attention are generally assumed to follow the biased competition principle [31]: Gains in processing [e.g., 32–34] for an attended location or feature are achieved by biasing neurons’ receptive fields in their favor, at the expense of unattended locations or features.

Without question, attention prior to or during encoding has high utility to behavior. However, the relevance of stimuli is not always obvious while they are still present – sometimes we need to prioritize information that has already been encoded in WM. For example, you may be looking around your apartment for your car keys and your phone simultaneously, holding templates of both in working memory as you scan your surroundings. Suddenly the phone starts ringing, so you prioritize finding the phone first to

get to it in time. This ability was already noted by William James in his endlessly cited definition of attention as the ‘taking possession by the mind [...] of one out of what seem several simultaneously possible objects or trains of thought’ [35] (p. 403-404). The ability flexibly to manipulate WM contents is also a hallmark of classic definitions of WM [1]. As with selective prioritization before and during encoding, prioritizing important items in WM during the retention interval has been shown to lead to a substantial memory boost [18–20]. Experimentally, prioritization within WM is generally induced by presenting a cue during the retention interval that directs focus to one of the items already held in mind. Cues can refresh a previously presented item [36,37], bring a subset of items currently in WM into the focus of attention (see Glossary) for immediate recall [38], or retroactively indicate that one item is most likely to be probed at the end of the delay interval. The latter is often referred to as a retrocue (as opposed to a precue presented before WM encoding, see [19], or a postcue presented together with the probe).

At first blush, the benefit of retrocuing seems paradoxical: memory is seemingly improved out of thin air. After all, the relevant information has already been stored in the brain, so how could providing an orienting cue possibly improve the strength of this information after the fact? Indeed, over ten years of investigation into the behavioral correlates and neural mechanisms of prioritization in WM have not yielded a conclusive explanation. Most proposals draw parallels between the effects of retrocuing and selective attention to external stimuli [9,13–15,39]. The same cognitive and neural mechanisms (selective attention) are deployed in each case, with the main difference being the substrate on which they operate, yielding a distinction between external and internal attention. In sum, these models emphasize that retrocuing benefits depend on a sustained bias of selective attention toward cued locations or features during a memory delay.

Overlap Between External and Internal Attention

At a basic level, the behavioral effects of retrocues indeed seem to be similar to the effects of external attention. Responses when probing cued items are faster and more accurate [18,19,40–42]. On invalid cue trials, responses are often slower and less accurate ([40,43–45], although invalidity costs are not consistent across studies, as discussed in the next section). When cue validity is manipulated, more reliable cues lead to a larger benefit [46–48].

At the neural level, similar top-down attention networks are engaged for internal and external attention shifts (such as the frontal eye fields and the superior parietal lobe, [49–59]). Neurophysiological markers of attention shifts, such as desynchronization of alpha-band oscillations in the hemisphere contralateral to where the cued item was presented, appear also to be roughly comparable between external attention [60] and retrocuing [59,61–64]. In parallel, retrocues seem to reduce load-dependent signals, such as the contralateral delay activity [65–67], as if they allowed the removal of uncued items from the memory store. This removal mechanism is reminiscent of the filtering of distractors during encoding [16]. In tandem with the top-down control signals, activity is also modulated in sensory brain regions corresponding to the cued location or feature [57,68–73], which likely contribute to WM representations or their manipulation through sensory recruitment [74–

79]. For example, when a visual stimulus category (e.g., faces) is cued, BOLD activity in the corresponding brain area (fusiform gyrus) increases [68,80,81]. This has been interpreted as increased processing of, or focused attention towards, the cued category. In many cases, this increase may, however, also reflect anticipation of a probe stimulus from that category at the end of the delay [57].

Open Questions for the Internal Attention Framework

For external information, attentional selection comes with a clear trade-off: attention to one object entails withdrawal from others. Selective attention is, to some extent, a zero-sum operation. By contrast, selecting an internal representation need not have this same constraint. Arguably the successful encoding and maintenance of individual items within working memory already entails a high degree of individuation and orthogonalization of their representations, thus decreasing the amount of potential interference among memoranda relative to what can occur during the encoding phase. Therefore, while selective biasing may still be in operation, the nature of the substrate is such that its consequences may be very different. In principle, at least, they could still be selected later, at low or no cost. From a functional perspective, it would be desirable to maintain memories for as long as they are potentially relevant to behavior, and only delete them once they are very unlikely to be useful. Therefore, while the biased-competition principle is a good starting point for proposing a mechanism of internally guided attention, its most basic component – selection via biased competition – may operate in crucially different ways on external vs. internal information.

The need for a distinction between external and internal attention has been highlighted before [13]: ‘Attention is not unitary’ (p. 76). We welcome such a careful differentiation between attention to perception and attention to working memory. We further propose that while both perception and working memory have limited capacity, the nature of the limit may differ considerably between the two. In perception, the challenge is to form cohesive and individuated item representations by bringing together their various attributes and separating these from competing sources of stimulation. In WM, the challenge is to select and use appropriate items to guide behavior.

Nominally, selection of one piece of information from among several in WM must occur in order to prioritize it (via a retrocue, for example). We propose, however, that this process differs at both the mechanistic and implementational levels from selection during perception. Our elaboration of the proposed mechanism follows in the next section.

Prioritizing Information in Working Memory

We propose that instead of invoking ‘internal attention,’ the prioritization of WM contents is better described as the attentional selection and, importantly, the reformatting of one out of several currently held memories to guide the next action. We speculate that this is a multi-step process. After a cue indicates the increased relevance of a particular item within working memory, the first step is to orient toward and select the cued item in the WM store. Orienting and selection can be thought of as the targeting of those neurons (for example, in

visual cortex) that are tuned to the location where the cued item was encountered and to the stimulus dimensions that are relevant to the memory (i.e., in a color WM task, activation might increase in color-sensitive visual areas such as V4). This allows for the effective grouping of all features belonging to the cued object [82], which in turn may reduce noise in the neural population representing it ([39], see also Box 2), potentially leading to increases in the precision of recalling cued items (e.g., [46,59,83,84]). Orienting attention in WM may be virtually identical, at the neural level, to the effects of preparatory orienting of attention for perception, and could explain the activation of canonical attention-control circuits after a retrocue, as summarized above. Critically, prioritization in WM allows for the immediate selection of memoranda. Selection of the cued representation can be thought of as an increased activation of the neurons coding for the relevant features of the cued object, possibly via reactivation of an ensemble that has encoded the item in its latent state (see Box 1 and [73]). Importantly, selection acts on one out of several objects in memory that have already been individuated and stored separately. This selection step therefore differs between prioritization in WM and in perception: In WM, the relevant information is already stored and can be selected immediately. During attention to external events, by contrast, selection cannot take place until the cued event occurs, and selection requires the identification of cued objects and their associated features among all perceptual input. While orienting and selection can be clearly delineated from one another in WM, for succinctness we will use ‘selection’ instead of ‘orienting and selection’ throughout the rest of this article.

Selection in memory appears to be a key element to successful prioritization in WM. Most accounts of retrocuing assume that this selection step is sufficient to account for the full range of experimental data. The essential novel aspect of our proposal is the speculation that the behavioral benefit additionally accrues in a further step. Following selection, the cued sensory representation can be transformed into a prospective, action-oriented representation, the better to influence behavior. By contrast, this step cannot take place in preparatory attention because the relevant information has not been presented yet. This transformation allows the current task set to become much more specific. For instance, in a typical visual WM change-detection experiment, the task set on a trial without a retrocue might reflect the following rule: ‘press button A if the probe stimulus matches the WM stimulus that was presented at the same location, and button B otherwise.’ Now imagine a retrocue trial, where one of the WM stimuli is cued (say it happens to be a green bar, see Fig 2). In our framework, the reformatted representation of the cued WM stimulus is now part of the task set. Therefore, the task set has become much more specific, and much simpler: ‘press button A if the probe stimulus is green, and button B otherwise.’ This process can be thought of as a form of cued task-set switching. When a task set is cued, responses are typically more efficient than when the task is not cued [85–88]. Therefore, the improved preparation for the task of responding accurately to the probe may, in part, contribute to the observed reaction-time and accuracy benefits. Crucially, once reformatting is complete, it is no longer necessary to sustain selective attention to the sensory representation that initially stored the cued information [89].

Multiple States of Representation in Working Memory

Our framework helps explain findings that are harder to reconcile with the prevailing account. For instance, in apparent contradiction of the sustained attention model, maintaining a constant attentional focus on cued representations is not necessary for retrocues to benefit behavior: After a retrocue has been fully processed, attention can be withdrawn from the cued item towards another task [90,91] or another WM representation [44] without impacting the retrocue benefit. These findings can be readily explained in our framework since, after prioritization and reformatting are complete, a sustained selective bias is no longer strictly necessary.

As we argue, retrocuing benefits arise in large part due to the prospective reformatting of the cued representation for use at the time of the probe. Reformatting appears to be a flexible process, meaning that other stored items could be prioritized at minimal cost. This is consistent with the finding that benefits can occur without costs to uncued items under some circumstances [18,46–48,92]. At least in principle, resource tradeoffs are not a necessity for benefits. Arguably, our framework is also consistent with the intuition that items held simultaneously in working memory are individually prioritized at different points in time as they become relevant to behavior [93]. For example, previously uncued items can be subsequently refocused by a second retrocue [44,45] or an internally generated change in expectation [94]. In each case, the retrocue effectively increases net WM capacity. In contrast to our proposal, a sustained internal attention mechanism that enhances the target representation and suppresses distractors should entail a trade-off in memory, and therefore cannot explain such findings as easily.

As mentioned above, invalid trials do seem to create costs in some studies. In our framework, such costs could still occur for a variety of reasons. Firstly, there are scenarios in which the initial selection step could create costs to uncued representations. For example, very high cue validity (e.g., when cues correctly predict the probe item on almost 100 percent of trials) may encourage a strategy of focusing all resources solely on the cued item by dropping uncued items from memory [83]. Such a strategy would be less successful when cues have lower validity, where it pays off to maintain uncued items in case of an invalid trial [46–48,95,96]. Additionally, high-validity cues might be used on a relatively higher proportion of trials, further increasing the retrocuing benefit (see Outstanding Questions).

The selection step might also account for effects of retrocuing on the precision of memory in some studies. Memory errors can be decomposed into errors due to Gaussian noise in the representation of the probed memory (i.e., its precision), errors due to forgetting, and errors due to misreporting the feature of an incorrect item [97,98]. The selection step, by strengthening the association between a cued location and the features of the object presented at that location, and by suppressing some of the residual interference from other items stored in WM, could lead to reductions in noise in the representation, leading to the occasional observation of small increases in memory precision [46,59,83,84]. However, the neural basis of this effect is unknown and difficult to fully explain, even within the sustained attention account. Finally, the selection step may also be used to select an ensemble of multiple items from working memory when multi-item cues are used [99]. In this case, it

appears that the entire set of cued items may be prioritized as an ensemble, rather than all cued items individually [100], possibly drawing on the visual system's representation of working memory contents at multiple, hierarchical organized spatial levels (from features bound to objects to ensembles of objects, [82,101,102]).

Furthermore, invalidity costs may also arise during the second step of reformatting the cued item. When a cued item is reformatted into an action-oriented format and a corresponding probe is anticipated, invalid trials may produce costs because of errors in probe anticipation and task preparation, instead of or in addition to any memory errors. This seems to fit the general pattern of the data: in single-cue studies, an invalid trial will occur when the incorrect task set is prepared, and therefore generates switch costs or response conflict [19,e.g., 40,62] which appear to have a particularly consistent influence on reaction time in addition to reducing memory accuracy somewhat [96]. This would be expected when an incorrect probe is expected while uncued items are still partially retained in memory. By contrast, studies employing a second cue in the delay that can redirect prioritization to a previously uncued item tend to find smaller or no costs on those trials [18,103–105]. Similarly, unanticipated task switches are known to incur behavioral costs [88,106,107]. In our framework, probing an uncued item amounts to an unanticipated task switch because the response must now be based on different information. Because the task-set representation is necessarily limited (since only one given set of rules should determine actions at any one time, especially if other rules would produce conflicting behavior), this could incur costs. Therefore, task-switching costs induced by setting up an inappropriate task set, over and above impaired memory alone, may explain performance decreases for uncued items. Importantly, this need not indicate a competition between the memory representations themselves. As a result, we expect cueing costs to be minimal when probe anticipation is controlled (for example, when a second cue cancels a prior retrocue).

In sum, our framework is consistent with a number of behavioral findings that are difficult to reconcile with a purely attentional account. Our framework, based on representational reformatting, relates to previous proposals [19,95] arguing that retrocues modulate the representation of a cued stimulus via attentional strengthening without necessarily requiring that other stimuli are deleted to provide more resources. Importantly, our framework makes several testable predictions about the nature of the cueing benefit. First, knowing the form of recall will affect the magnitude of the retrocue benefit because it will allow for more specific task set preparation. Second, several studies have shown that visual attention is drawn to items in the environment which match an item held in WM (in at least one of its features, e.g., location [12], color [108], etc.). It has been shown that this effect occurs only for items in the focus of attention [109]. Therefore, retrocued items should guide attention more than defocused items. However, this effect of attentional capture should also depend on the format in which they are about to be recalled. This seems to be the case [108], but has not been explicitly tested. For example, visual stimuli resembling a retrocued item should show increased attentional capture, compared to stimuli resembling unprioritized WM contents. However, our framework predicts that this capture effect should be larger if the WM task requires precise visual information, compared to when prioritized items can be maintained via a verbal strategy.

Neural Evidence for Multiple States in Working Memory

Neuroimaging studies support the existence of a second stage in prioritization of information in WM. Overall, these studies suggest that cueing a memory leads to reorganization of an output-oriented circuit which can then drive behavior faster and more accurately. These findings fall into two categories. The first is that output-related brain regions respond to retrocues and correlate with behavior, and the second is the additional activation of regions previously associated with task-set switching. We will discuss these sets of results in turn.

First, in addition to the well-documented activation of the top-down attention network (Fig 1), numerous studies have found additional activation in (primarily ventrolateral) prefrontal cortex [53–57,80,81] and striatum, which are less reliably related to external attention shifts. In one recent fMRI study involving retrocues and precues [110], retrocues led to correlations between the response strength of the caudate (as well as premotor cortex) and improvements in memory (as measured by reaction time). The authors [110] argued that this finding is consistent with the existence of output gating in working memory. The ‘output gate’ here can be thought of as a bottleneck forcing the selection of one item from all items currently held in working memory, so that it alone can guide the next action. The ‘output gate’ concept may be analogous to the behavior-guiding representational state proposed here. Other studies have also found striatal activation in response to retrocues on occasion [51], but this structure has generally been overlooked in discussions of the topic.

Second, the ventrolateral prefrontal cortex, stretching into the frontal operculum or anterior insula [111], is consistently activated in response to retrocues [49,53–57,59,80,81], and has been shown to be uniquely activated during retrocuing compared to external attention shifts (Fig 1, see [54]). The role of these areas is still somewhat unclear. Several studies have indicated that ventrolateral PFC is involved in the top-down access to and selection from sensory cortex of the cued information [80,81]. Consequently, disrupting activity in this area via transcranial magnetic stimulation reduces the benefits of retrocuing [81]. In addition, retrocuing tasks recruit task-switching-related brain circuits in lateral and medial prefrontal cortex [112]. This is consistent with our interpretation that transferring a retrocued item into the behavior-optimized state is akin to implementing a new task set.

These studies imply that additional prefrontal circuits are responsible for the top-down prioritization of items in working memory. However, fMRI lacks the temporal resolution to determine whether these additional areas and more traditional attention-related areas co-activate simultaneously, or whether they activate sequentially (as proposed in our framework). We predict that using selection to reactivate a memory representation is a transient process (step 1 - selection) that leads to a reconfiguration of the stimulus-response mappings of an action-selection network (step 2 - reformatting). After reformatting, sustained attention is no longer necessary (see Fig 2). One possibility is that cued information has been transferred to prefrontal cortex [2,41,113] via temporary synchronization of the cued region [114], and that after this process is completed, the attention-related modulation subsides. Studies investigating the neural basis of the focus of attention have found evidence that the lateral parietal cortex, rather than just prefrontal cortex, is critical for the deployment of this function [115–118]. The extent of network

interactions between these areas during prioritization remains to be fully investigated, but a recent study found evidence that frontal and parietal areas are both important for switching a working-memory representation into the focus of attention [119]. Behavioral data suggest that the benefits of retrocues emerge after 300 to 500 ms [120–122] – that is, the entire process of selection is completed within less than a second, making it difficult to use methods with low temporal resolution (such as BOLD fMRI) to settle the question of whether sustained attention to a cued feature is necessary for prioritized read-out.

Neuroimaging methods with the requisite temporal resolution, such as electroencephalography or magnetoencephalography, have shown that orienting and selection may be time-limited [59,63,123]. A recent MEG study [59] used lateralization of 10-Hz oscillations as a marker of internal attention shifts. The relative power of 10-Hz oscillations in sensory brain areas contralateral to where attention is directed, compared with power in ipsilateral areas, is a reliable indicator of covert attention shifts to external stimuli [60]. Similarly, cueing a location where a current WM item had previously appeared led to reliable lateralization. However, the lateralization of 10-Hz oscillations was transient after a retrocue, subsiding in less than a second. Given that the lateralization was only temporary, it seems unlikely that sustained attention or sustained active processing in a retinotopic representational format is necessary for retrocuing benefits to occur. Instead, this benefit could be conferred by a representational state change. Interestingly, additional activation in the insula that was specific to retrocues appeared only after the top-down attentional signal had peaked, supporting the idea of a two-step process of attentional selection, followed by representational reformatting [59]. The neural basis of representations after prioritization is still being investigated (see Box 2). A second study testing the efficacy of retrocues in elderly participants [63] confirmed the same temporary lateralization. Interestingly, the strength of the retrocuing benefit (the increase in accuracy compared to a neutral-cue) correlated negatively with the duration of the 10-Hz lateralization. Participants with the largest benefit showed the most transient lateralization: faster prioritization may therefore indicate a more accurate reconfiguration of the network. While this is consistent with our hypothesis that selection is only a temporary process, it seems to contradict the idea that sustained internal attention is crucial to improving behavior, because this should result in the exact opposite pattern: improved memory in those participants with more sustained 10-Hz lateralization.

Concluding Remarks

In summarizing the behavioral and neural literature on prioritization of working memory contents, we have argued that prioritizing an item arises in multiple stages and across multiple representational states. When the behavioral relevance of one out of several items in memory increases, top-down selection activates the neural subpopulation coding the cued information. Importantly, in a second step the cued information undergoes a transformation in its representational state, from a task-agnostic mnemonic representation to a task-specific representation that is best suited to guide behavior. This transformation may coincide with a transfer of information from sensory to action-guiding areas of the brain. Whether it also coincides with a temporary switch from latent to active representation (i.e., involving sustained spiking of memory-encoding neurons), however, is still an unresolved issue. With

recent advances in modeling memory-guided behavior and in multivariate analysis of high-dimensional neural recordings, we are confident that the predictions arising from our framework can be put to the test soon.

Acknowledgments

This article was funded by the Medical Research Council (to M.G.S.), the Wellcome Trust (A.C.N.: Senior Investigator Award 104571/Z/14/Z, and N.E.M.), University College, Oxford (N.E.M.), the James S. McDonnell Foundation (A.C.N.), and the National Institute for Health Research (NIHR) Oxford Biomedical Research Centre based at Oxford University Hospitals NHS Foundation Trust, Oxford University. The views expressed are those of the authors and not necessarily those of the NHS, the NIHR, or the Department of Health. We are grateful to Nahid Zokaei, Freek van Ede, and two anonymous reviewers for helpful comments on the manuscript, and to George Wallis for help with figures.

References

1. Baddeley AD, Hitch G. Working memory. *Psychology of learning and motivation*. 1974
2. Miller EK, Cohen JD. An Integrative Theory of Prefrontal Cortex Function. *Annu Rev Neurosci*. 2001; 24:167–202. [PubMed: 11283309]
3. Engle RW, et al. Working memory, short-term memory, and general fluid intelligence: A latent-variable approach. *Journal of Experimental Psychology: General*. 1999; 128:309–331. [PubMed: 10513398]
4. Smith PL, Ratcliff R. An integrated theory of attention and decision making in visual signal detection. *Psychological Review*. 2009; 116:283–317. [PubMed: 19348543]
5. Otto AR, et al. Working-memory capacity protects model-based learning from stress. *PNAS*. 2013; 110:20941–20946. [PubMed: 24324166]
6. Luck SJ, Vogel EK. Visual working memory capacity: from psychophysics and neurobiology to individual differences. *Trends in Cognitive Sciences*. 2013; 17:391–400. [PubMed: 23850263]
7. Phillips WA. On the distinction between sensory storage and short-term visual memory. *Perception & Psychophysics*. 1974; 16:283–290.
8. Cowan N. The magical number 4 in short-term memory: a reconsideration of mental storage capacity. *Behav Brain Sci*. 2000; 24:87–114.
9. Ma WJ, et al. Changing concepts of working memory. *Nature Neuroscience*. 2014; 17:347–356. [PubMed: 24569831]
10. Gazzaley A, Nobre AC. Top-down modulation: bridging selective attention and working memory. *Trends in Cognitive Sciences*. 2011; 16:129–135. [PubMed: 22209601]
11. Awh E, Jonides J. Overlapping mechanisms of attention and spatial working memory. *Trends in Cognitive Sciences*. 2001; 5:119–126. [PubMed: 11239812]
12. Awh E, et al. The role of spatial selective attention in working memory for locations: evidence from event-related potentials. *J Cog Neurosci*. 2000; 12:840–847.
13. Chun MM, et al. A Taxonomy of External and Internal Attention. *Annu Rev Psychol*. 2011; 62:73–101. [PubMed: 19575619]
14. Chun MM. Visual working memory as visual attention sustained internally over time. *Neuropsychologia*. 2011; 49:1407–1409. [PubMed: 21295047]
15. Kiyonaga A, Egner T. Working memory as internal attention: toward an integrative account of internal and external selection processes. *Psychonomic Bulletin & Review*. 2013; 20:228–242. [PubMed: 23233157]
16. Vogel EK, Machizawa MG. Neural activity predicts individual differences in visual working memory capacity. *Nature*. 2004; 428:748–751. [PubMed: 15085132]
17. Cusack R, et al. Encoding strategy and not visual working memory capacity correlates with intelligence. *Psychonomic Bulletin & Review*. 2009; 16:641–647. [PubMed: 19648446]
18. Landman R, et al. Large capacity storage of integrated objects before change blindness. *Vision Research*. 2003; 43:149–164. [PubMed: 12536137]

19. Griffin IC, Nobre AC. Orienting attention to locations in internal representations. *J Cog Neurosci*. 2003; 15:1176–1194.
20. Oberauer K, Hein L. Attention to Information in Working Memory. *Current Directions in Psychological Science*. 2012; 21:164–169.
21. Holmes J, et al. Working memory deficits can be overcome: Impacts of training and medication on working memory in children with ADHD. *Appl Cognit Psychol*. 2010; 24:827–836.
22. Gold JM, et al. Reduced capacity but spared precision and maintenance of working memory representations in schizophrenia. *Archives of General Psychiatry*. 2010; 67:570–577. [PubMed: 20530006]
23. Gazzaley A, et al. Top-down suppression deficit underlies working memory impairment in normal aging. *Nature Neuroscience*. 2005; 8:1298–1300. [PubMed: 16158065]
24. Reuter-Lorenz PA, Sylvester C-YC. The cognitive neuroscience of working memory and aging. *Cognitive Neuroscience of Aging: Linking Cognitive and Cerebral Aging*. 2005:186–217.
25. Pashler H. Processing stages in overlapping tasks: evidence for a central bottleneck. *Journal of Experimental Psychology: Human Perception and Performance*. 1984; 10:358–377. [PubMed: 6242412]
26. Raymond JE, et al. Temporary suppression of visual processing in an RSVP task: an attentional blink? *Journal of Experimental Psychology: Human Perception and Performance*. 1992; 18:849–860. [PubMed: 1500880]
27. Zanto TP, Gazzaley A. Neural Suppression of Irrelevant Information Underlies Optimal Working Memory Performance. *Journal of Neuroscience*. 2009; 29:3059–3066. [PubMed: 19279242]
28. Schmidt BK, et al. Voluntary and automatic attentional control of visual working memory. *Perc & Psychophys*. 2002; 64:754–763.
29. Murray AM, et al. Markers of preparatory attention predict visual short-term memory performance. *Neuropsychologia*. 2011; 49:1458–1465. [PubMed: 21335015]
30. Gazzaley A. Influence of early attentional modulation on working memory. *Neuropsychologia*. 2011; 49:1410–1424. [PubMed: 21184764]
31. Desimone R, Duncan J. Neural mechanisms of selective visual attention. *Annu Rev Neurosci*. 1995; 18:193–222. [PubMed: 7605061]
32. Hillyard SA, et al. Sensory gain control (amplification) as a mechanism of selective attention: electrophysiological and neuroimaging evidence. *Philosophical Transactions of the Royal Society B: Biological Sciences*. 1998; 353:1257–1270.
33. Nobre AC, et al. Modulation of human extrastriate visual processing by selective attention to colours and words. *Brain*. 1998; 121:1357–1368. [PubMed: 9679786]
34. Kastner S, Ungerleider LG. Mechanisms of visual attention in the human cortex. *Annu Rev Neurosci*. 2000; 23:315–341. [PubMed: 10845067]
35. James, W. *The Principles of Psychology*. Henry Holt; 1890.
36. Johnson MK, et al. Second Thoughts versus Second Looks: An Age-Related Deficit in Reflectively Refreshing Just-Activated Information. *Psychological Science*. 2002; 13:64–67. [PubMed: 11892780]
37. Chun MM, Johnson MK. Memory: Enduring Traces of Perceptual and Reflective Attention. *Neuron*. 2011; 72:520–535. [PubMed: 22099456]
38. Oberauer K. Access to information in working memory: Exploring the focus of attention. *Journal of Experimental Psychology: Learning, Memory, and Cognition*. 2002; 28:411–421.
39. Bays PM. Spikes not slots: noise in neural populations limits working memory. *Trends in Cognitive Sciences*. 2015; 19:431–438. [PubMed: 26160026]
40. Astle DE, et al. Orienting attention to locations in mental representations. *Atten Percept Psychophys*. 2012; 74:146–162. [PubMed: 21972046]
41. Sligte IG, et al. Detailed sensory memory, sloppy working memory. *Front Psychology*. 2010; 1:175.
42. Makovski T, et al. Orienting attention in visual working memory reduces interference from memory probes. *Journal of Experimental Psychology: Learning, Memory, and Cognition*. 2008; 34:369–380.

43. Gözenman F, et al. Invalid retro-cues can eliminate the retrocue benefit: Evidence for a hybridized account. *Journal of Experimental Psychology: Human Perception and Performance*. 2014; 40:1748–1754. [PubMed: 25045904]
44. Rerko L, et al. Retro-cue benefits in working memory without sustained focal attention. *Memory & Cognition*. 2014; 42:712–728. [PubMed: 24442551]
45. van Moorselaar D, et al. Forgotten but not gone: Retro-cue costs and benefits in a double-cueing paradigm suggest multiple states in visual short-term memory. *Journal of Experimental Psychology: Learning, Memory, and Cognition*. 2015; 41:1755–1763.
46. Gunseli E, et al. The reliability of retro-cues determines the fate of noncued visual working memory representations. *Psychonomic Bulletin & Review*. 2015; 22:1334–1341. [PubMed: 25563713]
47. Berryhill ME, et al. Shifting attention among working memory representations: Testing cue type, awareness, and strategic control. *The Quarterly Journal of Experimental Psychology*. 2012; 65:426–438. [PubMed: 21846267]
48. Shimi A, et al. Orienting attention within visual short-term memory: development and mechanisms. *Child Dev*. 2014; 85:578–592. [PubMed: 23937596]
49. Nobre AC, et al. Orienting Attention to Locations in Perceptual Versus Mental Representations. *J Cog Neurosci*. 2004; 16:363–373.
50. Bledowski C, et al. What “Works” in Working Memory? Separate Systems for Selection and Updating of Critical Information. *Journal of Neuroscience*. 2009; 29:13735–13741. [PubMed: 19864586]
51. Tamber-Rosenau BJ, et al. Cortical Mechanisms of Cognitive Control for Shifting Attention in Vision and Working Memory. *J Cog Neurosci*. 2011; 23:2905–2919.
52. Nee DE, et al. A Meta-analysis of Executive Components of Working Memory. *Cerebral Cortex*. 2013; 23:264–282. [PubMed: 22314046]
53. Nee DE, Jonides J. Neural correlates of access to short-term memory. *PNAS*. 2008; 105:14228–14233. [PubMed: 18757724]
54. Nee DE, Jonides J. Common and distinct neural correlates of perceptual and memorial selection. *NeuroImage*. 2009; 45:963–975. [PubMed: 19280708]
55. Lepsien J, Nobre AC. Cognitive control of attention in the human brain: Insights from orienting attention to mental representations. *Brain Research*. 2006; 1105:20–31. [PubMed: 16729979]
56. Lepsien J, et al. Directing spatial attention in mental representations: Interactions between attentional orienting and working-memory load. *NeuroImage*. 2005; 26:733–743. [PubMed: 15955482]
57. Lepsien J, et al. Modulation of working-memory maintenance by directed attention. *Neuropsychologia*. 2011; 49:1569–1577. [PubMed: 21420421]
58. Yeh Y-Y, et al. The neural correlates of attention orienting in visuospatial working memory for detecting feature and conjunction changes. *Brain Research*. 2007; 1130:146–157. [PubMed: 17173876]
59. Wallis G, et al. Frontoparietal and Cingulo-opercular Networks Play Dissociable Roles in Control of Working Memory. *J Cog Neurosci*. 2015; 15:1–16.
60. Worden MS, et al. Anticipatory biasing of visuospatial attention indexed by retinotopically specific-band electroencephalography increases over occipital cortex. *J Neurosci*. 2000; 20:1–6. [PubMed: 10627575]
61. Poch C, et al. Modulation of alpha and gamma oscillations related to retrospectively orienting attention within working memory. *Eur J Neurosci*. 2014; 40:2399–2405. [PubMed: 24750388]
62. Myers NE, et al. Temporal Dynamics of Attention during Encoding versus Maintenance of Working Memory: Complementary Views from Event-related Potentials and Alpha-band Oscillations. *J Cog Neurosci*. 2015; 27:492–508.
63. Mok RM, et al. Behavioral and Neural Markers of Flexible Attention over Working Memory in Aging. *Cereb Cortex*. 2016; 26:1831–1842. [PubMed: 26865653]
64. van Ede F, et al. Temporal expectations guide dynamic prioritization in visual working memory through attenuated alpha oscillations. *J Neurosci*. 2016; doi: 10.1523/JNEUROSCI.2272-16.2016

65. Kuo B-C, et al. Attention Modulates Maintenance of Representations in Visual Short-term Memory. *J Cog Neurosci*. 2012; 24:51–60.
66. Nobre AC. Spatial attention can bias search in visual short-term memory. *Front Hum Neurosci*. 2008; 1
67. Williams M, Woodman GF. Directed forgetting and directed remembering in visual working memory. *Journal of Experimental Psychology: Learning, Memory, and Cognition*. 2012; 38:1206–1220.
68. Lepsien J, Nobre AC. Attentional Modulation of Object Representations in Working Memory. *Cerebral Cortex*. 2007; 17:2072–2083. [PubMed: 17099066]
69. Kuo B-C, et al. Attention Biases Visual Activity in Visual STM. *J Cog Neurosci*. 2014; 26:1377–1389.
70. Lewis-Peacock JA, Postle BR. Decoding the internal focus of attention. *Neuropsychologia*. 2012; 50:470–478. [PubMed: 22108440]
71. Lewis-Peacock JA, et al. Neural Evidence for the Flexible Control of Mental Representations. *Cerebral Cortex*. 2014; doi: 10.1093/cercor/bhu130
72. LaRocque JJ, et al. Decoding attended information in short-term memory: an EEG study. *J Cog Neurosci*. 2013; 25:127–142.
73. Sprague TC, et al. Restoring Latent Visual Working Memory Representations in Human Cortex. *Neuron*. 2016; 91:694–707. [PubMed: 27497224]
74. Pasternak T, Greenlee MW. Working memory in primate sensory systems. *Nat Rev Neurosci*. 2005; 6:97–107. [PubMed: 15654324]
75. Scolari M, Serences JT. Basing Perceptual Decisions on the Most Informative Sensory Neurons. *Journal of Neurophysiology*. 2010
76. Riggall AC, Postle BR. The relationship between working memory storage and elevated activity as measured with functional magnetic resonance imaging. *Journal of Neuroscience*. 2012; 32:12990–12998. [PubMed: 22993416]
77. Christophel TB, et al. The Distributed Nature of Working Memory. *Trends in Cognitive Sciences*. 2017; 21:111–124. [PubMed: 28063661]
78. Lee S-H, et al. Goal-dependent dissociation of visual and prefrontal cortices during working memory. *Nat Neuro*. 2013; 16:997–999.
79. Ester EF, et al. Parietal and Frontal Cortex Encode Stimulus-Specific Mnemonic Representations during Visual Working Memory. *Neuron*. 2015; 87:893–905. [PubMed: 26257053]
80. Nelissen N, et al. Frontal and parietal cortical interactions with distributed visual representations during selective attention and action selection. *Journal of Neuroscience*. 2013; 33:16443–16458. [PubMed: 24133250]
81. Higo T, et al. Distributed and causal influence of frontal operculum in task control. *Proceedings of the National Academy of Sciences*. 2011; 108:4230–4235.
82. Brady TF, et al. A review of visual memory capacity: Beyond individual items and toward structured representations. *Journal of Vision*. 2011; 11:4–4.
83. Williams M, et al. The benefit of forgetting. *Psychonomic Bulletin & Review*. 2013; 20:348–355. [PubMed: 23208769]
84. Pertzov Y, et al. Rapid Forgetting Results From Competition Over Time Between Items in Visual Working Memory. *Journal of Experimental Psychology: Learning, Memory, and Cognition*. 2016; doi: 10.1037/xlm0000328
85. Meiran N. Reconfiguration of processing mode prior to task performance. *Journal of Experimental Psychology: Learning, Memory, and Cognition*. 1996; 22:1423–1442.
86. Rogers RD, Monsell S. Costs of a predictable switch between simple cognitive tasks. *Journal of Experimental Psychology: General*. 1995; 124:207–231.
87. Monsell S. Task switching. *Trends in Cognitive Sciences*. 2003; 7:134–140. [PubMed: 12639695]
88. Allport, A, , et al. Shifting intentional set: Exploring the dynamic control of tasksAttention and Performance. Umiltà, CA, editor. Vol. 15. 1994. 412–452. *Attention and performance XV: Conscious and ...*

89. Muhle-Karbe PS, et al. Neural Coding for Instruction-Based Task Sets in Human Frontoparietal and Visual Cortex. *Cereb Cortex*. 2016; doi: 10.1093/cercor/bhw032
90. Hollingworth A, Maxcey-Richard AM. Selective maintenance in visual working memory does not require sustained visual attention. *Journal of Experimental Psychology: Human Perception and Performance*. 2013; 39:1047–1058. [PubMed: 23067118]
91. Makovski T, Pertzov Y. Attention and memory protection: Interactions between retrospective attention cueing and interference. *Q J Exp Psychol (Hove)*. 2015; 68:1735–1743. [PubMed: 25980784]
92. Souza AS, et al. Refreshing memory traces: thinking of an item improves retrieval from visual working memory. *Annals of the New York Academy of Sciences*. 2015; 1339:20–31. [PubMed: 25557544]
93. Duncan J. The Structure of Cognition: Attentional Episodes in Mind and Brain. *Neuron*. 2013; 80:35–50. [PubMed: 24094101]
94. van Ede F, et al. Temporal expectations guide dynamic prioritization in visual working memory through attenuated alpha oscillations. *Journal of Neuroscience*.
95. Souza AS, Oberauer K. In search of the focus of attention in working memory: 13 years of the retro-cue effect. *Atten Percept Psychophys*. 2016; doi: 10.3758/s13414-016-1108-5
96. Gressmann M, Janczyk M. The (Un)Clear Effects of Invalid Retro-Cues. *Front Psychology*. 2016; 7:244.
97. Zhang W, Luck SJ. Discrete fixed-resolution representations in visual working memory. *Nature*. 2008; 453:233–235. [PubMed: 18385672]
98. Bays PM, et al. The precision of visual working memory is set by allocation of a shared resource. *Journal of Vision*. 2009; 9:7–7.
99. Matsukura M, et al. Attention effects during visual short-term memory maintenance: protection or prioritization? *Perception & Psychophysics*. 2007; 69:1422–1434. [PubMed: 18078232]
100. Matsukura M, Vecera SP. Selection of multiple cued items is possible during visual short-term memory maintenance. *Atten Percept Psychophys*. 2015; doi: 10.3758/s13414-015-0836-2
101. Brady TF, Alvarez GA. Hierarchical Encoding in Visual Working Memory: Ensemble Statistics Bias Memory for Individual Items. *Psychological Science*. 2011; 22:384–392. [PubMed: 21296808]
102. Brady TF, et al. Compression in visual working memory: Using statistical regularities to form more efficient memory representations. *Journal of Experimental Psychology: General*. 2009; 138:487–502. [PubMed: 19883132]
103. Li Q, Saiki J. The effects of sequential attention shifts within visual working memory. *Front Psychology*. 2014; 5:965.
104. Rerko L, Oberauer K. Focused, unfocused, and defocused information in working memory. *Journal of Experimental Psychology: Learning, Memory, and Cognition*. 2013; 39:1075–1096.
105. Souza AS, Oberauer K. Time-based forgetting in visual working memory reflects temporal distinctiveness, not decay. *Psychonomic Bulletin & Review*. 2015; 22:156–162. [PubMed: 24825306]
106. Rushworth MFS, et al. Components of Switching Intentional Set. *J Cog Neurosci*. 2002; 14:1139–1150.
107. Miniussi C, et al. Modulation of brain activity by selective task sets observed using event-related potentials. *Neuropsychologia*. 2005; 43:1514–1528. [PubMed: 15989941]
108. Olivers CNL, et al. Feature-based memory-driven attentional capture: visual working memory content affects visual attention. *Journal of Experimental Psychology: Human Perception and Performance*. 2006; 32:1243–1265. [PubMed: 17002535]
109. Olivers CNL, et al. Different states in visual working memory: when it guides attention and when it does not. *Trends in Cognitive Sciences*. 2011; 15:327–334. [PubMed: 21665518]
110. Chatham CH, et al. Corticostriatal Output Gating during Selection from Working Memory. *Neuron*. 2014; 81:930–942. [PubMed: 24559680]

111. Petrides M, Pandya DN. Comparative cytoarchitectonic analysis of the human and the macaque ventrolateral prefrontal cortex and corticocortical connection patterns in the monkey. *Eur J Neurosci.* 2002; 16:291–310. [PubMed: 12169111]
112. Richter, FR, Yeung, N. *Neuroimaging Studies of Task Switching*. Task Switching and Cognitive Control. Grange, J, Houghton, G, editors. Oxford University Press; 2014. 1–55.
113. Buschman TJ, Miller EK. Goal-direction and top-down control. *Philos Trans R Soc Lond, B, Biol Sci.* 2014; 369:20130471–20130471. [PubMed: 25267814]
114. Liebe S, et al. Theta coupling between V4 and prefrontal cortex predicts visual short-term memory performance. *Nat Neuro.* 2012; doi: 10.1038/nn.3038
115. Öztekin I, Cowan N. Editorial: Representational states in memory: where do we stand? *Front Hum Neurosci.* 2015; 9:453. [PubMed: 26347637]
116. Nee DE, Jonides J. Trisecting representational states in short-term memory. *Front Hum Neurosci.* 2013; doi: 10.3389/fnhum.2013.00796/abstract
117. Nee DE, Jonides J. Neural evidence for a 3-state model of visual short-term memory. *NeuroImage.* 2013; 74:1–11. [PubMed: 23435212]
118. Nee DE, Jonides J. Dissociable contributions of prefrontal cortex and the hippocampus to short-term memory: evidence for a 3-state model of memory. *NeuroImage.* 2011; 54:1540–1548. [PubMed: 20832478]
119. Nee DE, Brown JW. Dissociable frontal-striatal and frontal-parietal networks involved in updating hierarchical contexts in working memory. *Cereb Cortex.* 2013; 23:2146–2158. [PubMed: 22798339]
120. Tanoue RT, Berryhill ME. The mental wormhole: Internal attention shifts without regard for distance. *Atten Percept Psychophys.* 2012; 74:1199–1215. [PubMed: 22549808]
121. Souza AS, et al. Focused attention improves working memory: implications for flexible-resource and discrete-capacity models. *Atten Percept Psychophys.* 2014; 76:2080–2102. [PubMed: 24874258]
122. van Moorselaar D, et al. The time course of protecting a visual memory representation from perceptual interference. *Front Hum Neurosci.* 2014; 8:1053. [PubMed: 25628555]
123. Spitzer B, Blankenburg F. Stimulus-dependent EEG activity reflects internal updating of tactile working memory in humans. *PNAS.* 2011; 108:8444–8449. [PubMed: 21536865]
124. Crone EA. Neural Evidence for Dissociable Components of Task-switching. *Cerebral Cortex.* 2005; 16:475–486. [PubMed: 16000652]
125. Buschman TJ, et al. Synchronous Oscillatory Neural Ensembles for Rules in the Prefrontal Cortex. *Neuron.* 2012; 76:838–846. [PubMed: 23177967]
126. Stokes MG. “Activity-silent” working memory in prefrontal cortex: a dynamic coding framework. *Trends in Cognitive Sciences.* 2015; 19:394–405. [PubMed: 26051384]
127. Funahashi S, et al. Prefrontal neuronal activity in rhesus monkeys performing a delayed anti-saccade task. *Nature.* 1993; 365:753–756. [PubMed: 8413653]
128. Mongillo G, et al. Synaptic Theory of Working Memory. *Science.* 2008; 319:1543–1546. [PubMed: 18339943]
129. Kaufman MT, et al. Cortical activity in the null space: permitting preparation without movement. *Nature Neuroscience.* 2014; 17:440–448. [PubMed: 24487233]
130. Bhandari A, Duncan J. Goal neglect and knowledge chunking in the construction of novel behaviour. *Cognition.* 2014; 130:11–30. [PubMed: 24141034]
131. Ericsson KA, Kintsch W. Long-term working memory. *Psychological Review.* 1995; 102:211–245. [PubMed: 7740089]
132. Abbott LF, Regehr WG. Synaptic computation. *Nature.* 2004; 431:796–803. [PubMed: 15483601]
133. Citri A, Malenka RC. Synaptic plasticity: multiple forms, functions, and mechanisms. *Neuropsychopharmacology.* 2008; 33:18–41. [PubMed: 17728696]
134. Rose NS, et al. Reactivation of latent working memories with transcranial magnetic stimulation. *Science.* 2016; 354:1136–1139. [PubMed: 27934762]
135. Postle BR. The cognitive neuroscience of visual short-term memory. *Current Opinion in Behavioral Sciences.* 2015; 1:40–46. [PubMed: 26516631]

136. Lewis-Peacock JA, et al. Neural Evidence for a Distinction between Short-term Memory and the Focus of Attention. *J Cog Neurosci*. 2012; 24:61–79.

Box 1**Latent Versus Active Neural States Supporting Working Memory**

Our framework predicts two functional states in WM – maintenance of information without a specific action plan, and of a prioritized item in an action-oriented format. We propose that the latter depends on flexible changes in the tuning of an action-oriented network that includes lateral prefrontal cortex [124,125]. WM has been proposed to rely on changes in the underlying state of a neural ensemble [126], permitting latent storage without requiring sustained activity [127]. The latent state change could rely on many physiological mechanisms. What they have in common is that they do not depend on an unbroken chain of sustained spiking. One candidate is short-term synaptic plasticity: the connectivity within a neural ensemble changes such that its response to subsequent input reflects the WM content [128]. Alternatively, memory-specific ensembles could emerge by synchronizing to a common oscillatory rhythm [125].

The concept of latent storage has been applied to maintenance of WM items per se, but latent storage may also be of particular relevance for representing prioritized WM items. Prioritization may lead to the task-specific transfer of a latent code stored in visual cortex to a lateral prefrontal network, possibly through temporary changes in the synaptic weights in PFC. Alternatively, unprioritized items could already be represented in PFC activity, in patterns that do not drive downstream motor regions [129] until they are prioritized.

After reconfiguration, sustained firing is unnecessary, so both prioritized and unprioritized representations are maintained using latent storage. Crucially, only the former may be represented in a format that is optimized for behavior. By contrast, unprioritized representations might be stored without immediately influencing behavior [130].

Latent storage of WM in changed connection weights invokes comparisons to long-term memory (LTM), which may operate along similar lines.

Prominent WM theories propose activated LTM as the basis of WM [8,20,131]. While latent WM and LTM storage may both depend on changing synaptic weights instead of persistent firing, the specifics of how and for how long synaptic weights are changed may differ considerably between WM and LTM [132,133]. For example, any synaptic weight changes subserving WM must be short-lived to avoid interference from traces of recent WM contents. This is not a constraint for LTM. Further, WM prioritization may reconfigure the action-oriented network in PFC so it can immediately produce a context-appropriate response to the probe, without first needing to recall information from LTM (as proposed by some WM models, e.g. [131]).

Box 2**Working Memory as Internal Attention**

Undoubtedly we have gained much from drawing parallels between internal and external attention. Much like selective attention towards perceptual representations is thought to bias competitive processing in favor of one representation over another, internal attention is argued to bias processing towards one mnemonic representation over others in a shared memory store [9,14,15,39]. The shift of resources improves retention of the cued item or the behavior guided by it. An influential review [13] argues that attention shares common principles across the substrates it acts on. What is shared are the purpose of attention (overcoming limited capacity via selection), and its consequence (modulation of the selected information). The core of this process is that ‘multiple stimuli [...] compete for selection, and the goal of attention is to bias competition in favor of a target object’ (p. 75). Therefore, ‘selecting a memory from competing memories should be viewed as an attentional operation. The cost is that unattended information may be missed.’

An extension of the internal attention account is that prioritizing a WM representation may equate to transforming it from a latent to an activated neural state [15,73,109,134,135], rather than transforming it to an output-oriented representation (Box 1). Thus, active versus latent storage corresponds at the neural level with attended versus unattended WM states at the cognitive level. In a recent study [73], retrocues improved decoding of the cued item in retinotopic visual areas. The authors argued that this finding is consistent with the reactivation of a latent code in sensory brain regions (the activated population permitting improved decoding, as noted elsewhere, [136]). However, the sluggishness of the fMRI signal may have precluded them from testing whether the reactivation reflects a temporary process. Most importantly, it is unclear how moving from latent to active representation alone could account for improved recall accuracy without simultaneously decreasing memory for uncued representations. For example, increasing the activity of a neural ensemble might give it greater influence over a downstream readout population, compared to a competing ensemble encoding an uncued representation (i.e., biased competition, [31]). Alternatively, activation might suppress activity of competing ensembles via lateral inhibition [39]. Either way, the increased activity confers a benefit only by virtue of its suppression of competitors.

Box 3**Outstanding Questions**

1. What is the representational format and neural substrate of a prioritized WM item? How does this format relate to the representation of task sets or task rules?
2. What is the relationship between cue validity and the size of the retrocueing effect?
3. Are similar prioritization mechanisms also important in preparing for more classic forms of WM manipulation (e.g., mental arithmetic)?
4. In reality, we experience a continuous stream of thoughts passing through WM. How do we extend the concept of flexible prioritization to continuous, temporally extended cognition?
5. How do we switch between an internal and an external focus?
6. Do long-term memory and working memory share selection and prioritization mechanisms? Furthermore, how does long-term memory influence the interplay between perception and WM?
7. Retrieval from long-term memory can induce forgetting of associated memories. Does a similar phenomenon exist in WM?
8. Can we dissociate the short-term representation of task goals or rules from the representation of other kinds of content in WM? Does the neural dissociability of goals and content depend on the task context?
9. WM is sustained by several representational states. Which of these corresponds to the traditional notion of the attention-guiding template? What other sources of attentional guidance exist, and what can the fractionation of WM tell us about the fractionation of the control of attention?

Trends Box

- Recent research has uncovered our remarkable flexibility in prioritizing information in working memory (WM), refining the concept of multiple representational states in WM.
- Neuroimaging studies have investigated the networks controlling prioritization in WM.
- Prioritization activates prefrontal and parietal brain areas associated with the deployment of visual attention, suggesting a parallel between attention to external stimuli and attention to memory contents ('internal attention').
- However, additional prefrontal areas are specific to WM prioritization. We propose that they reflect recruitment of high-priority information for the next action. What can this tell us about the neural basis of different representational states in WM? We speculate that prioritized information is reflected in the task-specific tuning of a neural network important for action selection and preparation.

Glossary

Focus of Attention: specialized state within working memory. As opposed to items that are merely maintained, the single item in the focus of attention [20] is selected and elevated to a separate representational state so that it can be updated, manipulated, or recalled. Representations in the focus of attention are recalled more quickly and with greater accuracy than other WM representations.

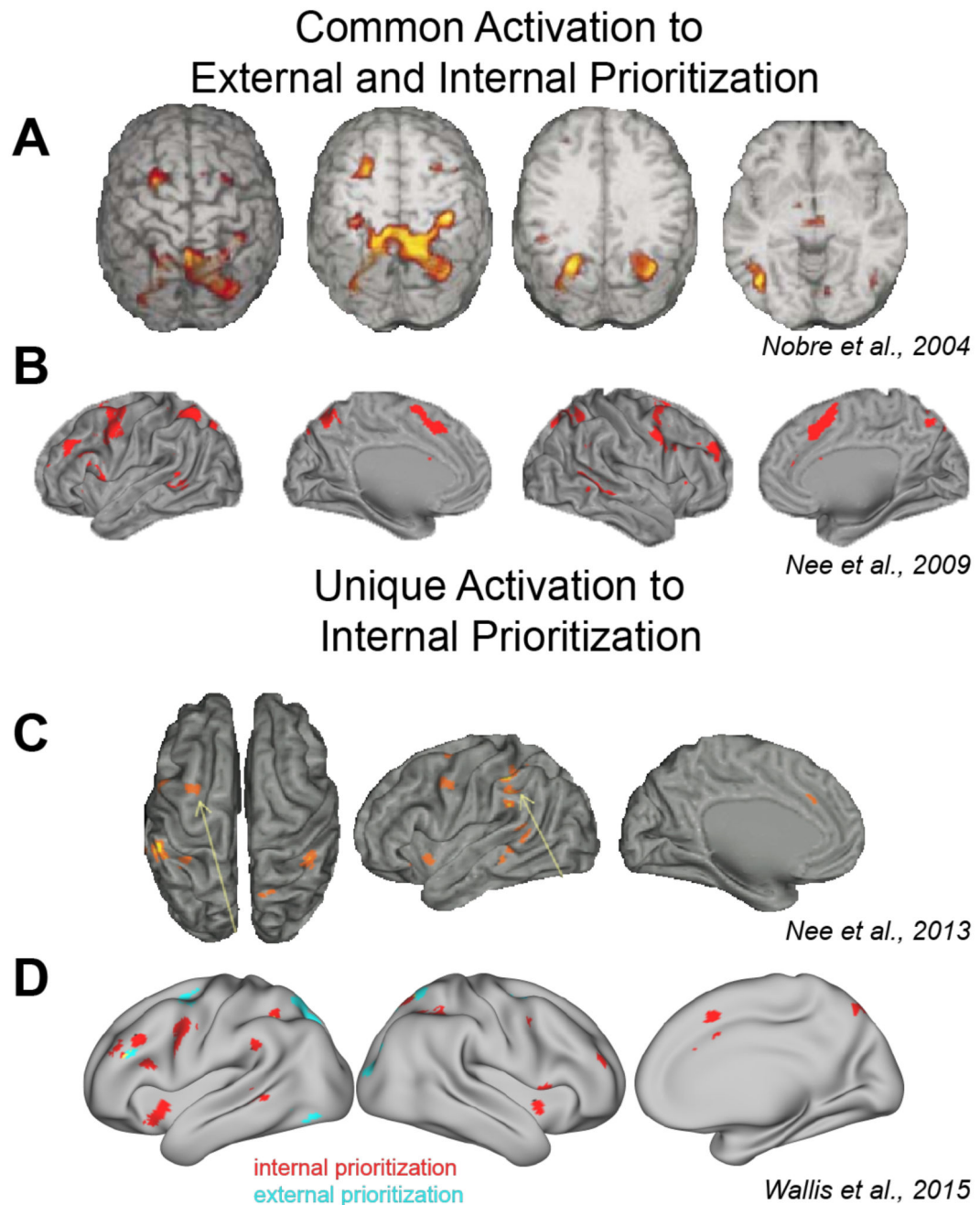
Internal Attention: goal-directed selection of information that is not currently presented in the environment, such as long- or short-term memory or goals. Internal attention is thought to draw on the same selection mechanism that is deployed to attend to information arriving from the environment.

Latent Storage: Proposed neurophysiological mechanism for the neural storage of WM memoranda by reconfiguring the state of a memory network through short-term changes in its pattern of connections. After reconfiguration, persistent spiking is no longer necessary because the memory is stored in a latent state, for example in temporarily changed synaptic weights.

Output Gating: Some computational models of WM emphasize the importance of an input gate that determines which pieces of sensory information are allowed into the limited-capacity WM store. Similarly, more recent computational models propose a second gate determining which of the items that are stored in WM are permitted to drive behavior, or 'output'. 'Output-gating' an item could correspond to moving it into the focus of attention, although the exact relationship is unclear.

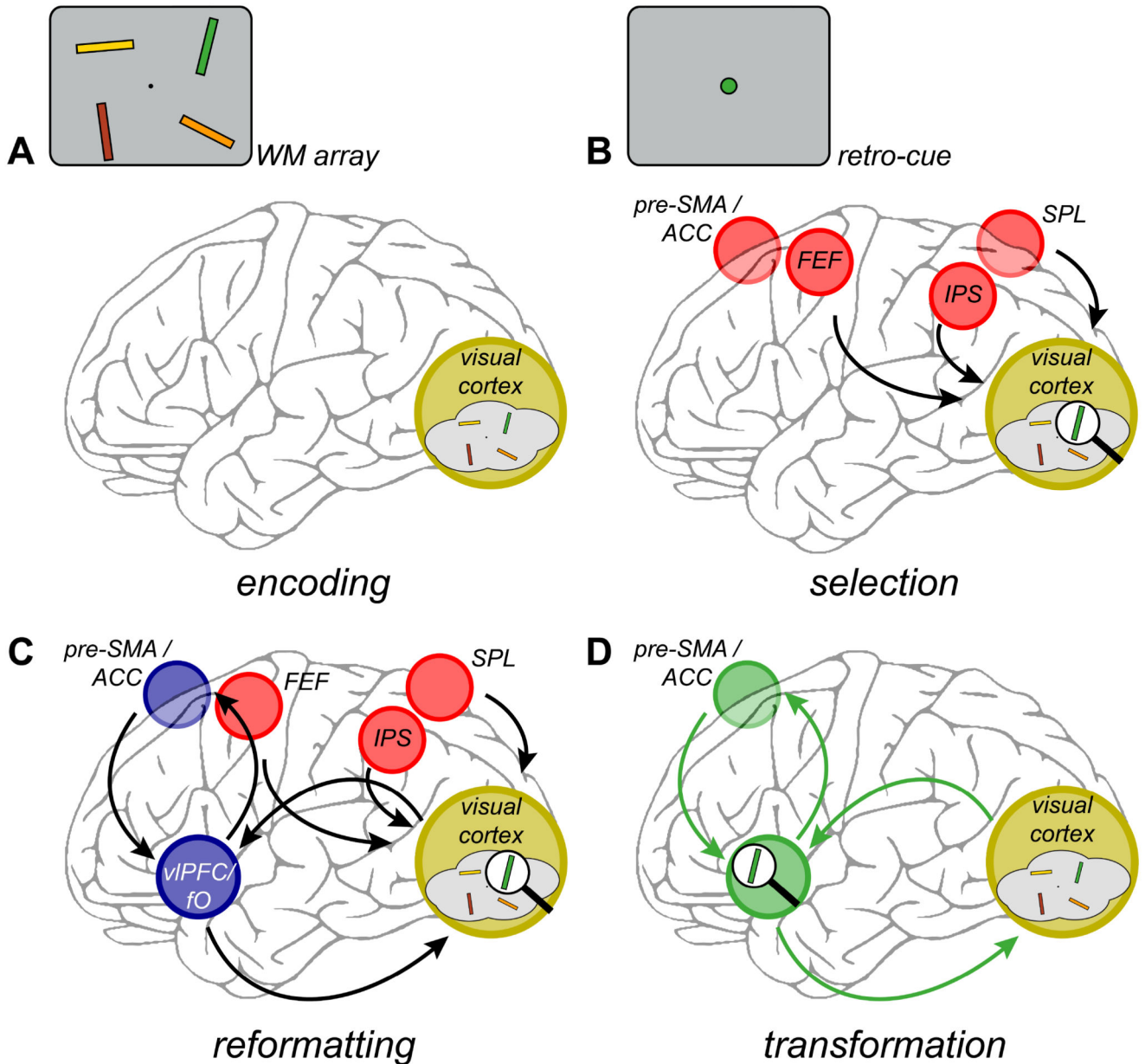
Retrocue: A cue presented retrospectively, during the retention interval of a working memory task, indicating that a subset of all items already held in memory is most relevant to behavior, for example because it is most likely to be probed.

Task Switching: Switching tasks (rule-guided responses to a limited set of stimuli) incurs costs in terms of slower reaction times and increased error rates. Switch costs occur because of the sudden need to reconfigure a task set in response-guiding brain networks. Cueing a task switch in advance reduces but does not eliminate switching costs.

**Figure 1.**

Shared and unique networks for attentional selection and prioritization in working memory. A. Spatial cues directing attention to external stimuli or to contents in working memory both activate a network spanning frontal eye fields, the pre-supplementary motor cortex and anterior cingulate, the intraparietal sulcus, and the superior parietal lobule [49]. B. This overlap has been confirmed in multiple neuroimaging studies [54] and a meta-analysis [59]. C. Additional areas respond only to prioritization within working memory, including ventrolateral prefrontal cortex stretching into the frontal operculum and anterior insula. In

the pre-SMA and ACC, activation is either stronger than during external attention shifts, or additional subregions are recruited [52]. D. A recent meta-analysis [59] found several mostly prefrontal areas responding to retrocues (internal prioritization, red) but not to precues (external prioritization, blue).

**Figure 2.**

Proposed sequence of prioritization in working memory. A. During encoding, sensory brain areas (in yellow) are recruited and modified to reflect the relevant features of a WM array. In the task illustrated here, sensory brain areas are the most likely substrate of memory maintenance because it requires memory for fine sensory details. However, the same mechanism could act on other brain areas if WM contents are stored elsewhere. B. A cue indicating that one WM item is of particular relevance (the green oriented bar, in this case) leads to orientation toward, and selection of the relevant representation. This operation recruits the top-down attention network that is also involved in external attention shifts (red circles). C. In a second step, the identified information is prioritized. This step recruits a

prefrontal network (in blue) comprising the anterior insula or frontal operculum and ventrolateral PFC (vIPFC/fO), and possibly pre-supplementary motor area or anterior cingulate cortex (pre-SMA/ACC). D. Finally, the selected representation is now reformatted to bring the network into an optimal state to respond to the recall demands of the probe stimulus.